

# Foundations of Machine Learning

Mohammad Emtiyaz Khan  
RIKEN-AIP, Tokyo  
OIST, Okinawa

<http://emtiyaz.github.io>

[emtiyaz.khan@riken.jp](mailto:emtiyaz.khan@riken.jp)

May 12, 2022



©Mohammad Emtiyaz Khan 2022

Hello!

# Goals

Understand (some) fundamentals of Machine learning<sup>1</sup>.

**Part I : Understand the **basic set-up** to analyze data under a machine-learning framework.**

1. Before Machine Learning. *Problem*

2. ML Problem: Regression. *Data*

3. Model: Linear Regression. *Model*

4. Cost Function: MSE.

5. Algorithm 1: Gradient Descent.

6. Algorithm 2: Least Squares.

*algorithm*

*"Foundation"*

**Part II : Understand what can go wrong when learning from data and how to correct it.**

6. Challenge: Overfitting.

7. Solutions: Regularization.

8. Bias-Variance Decomposition. *"Trade-off"*

9. Recent Advances.

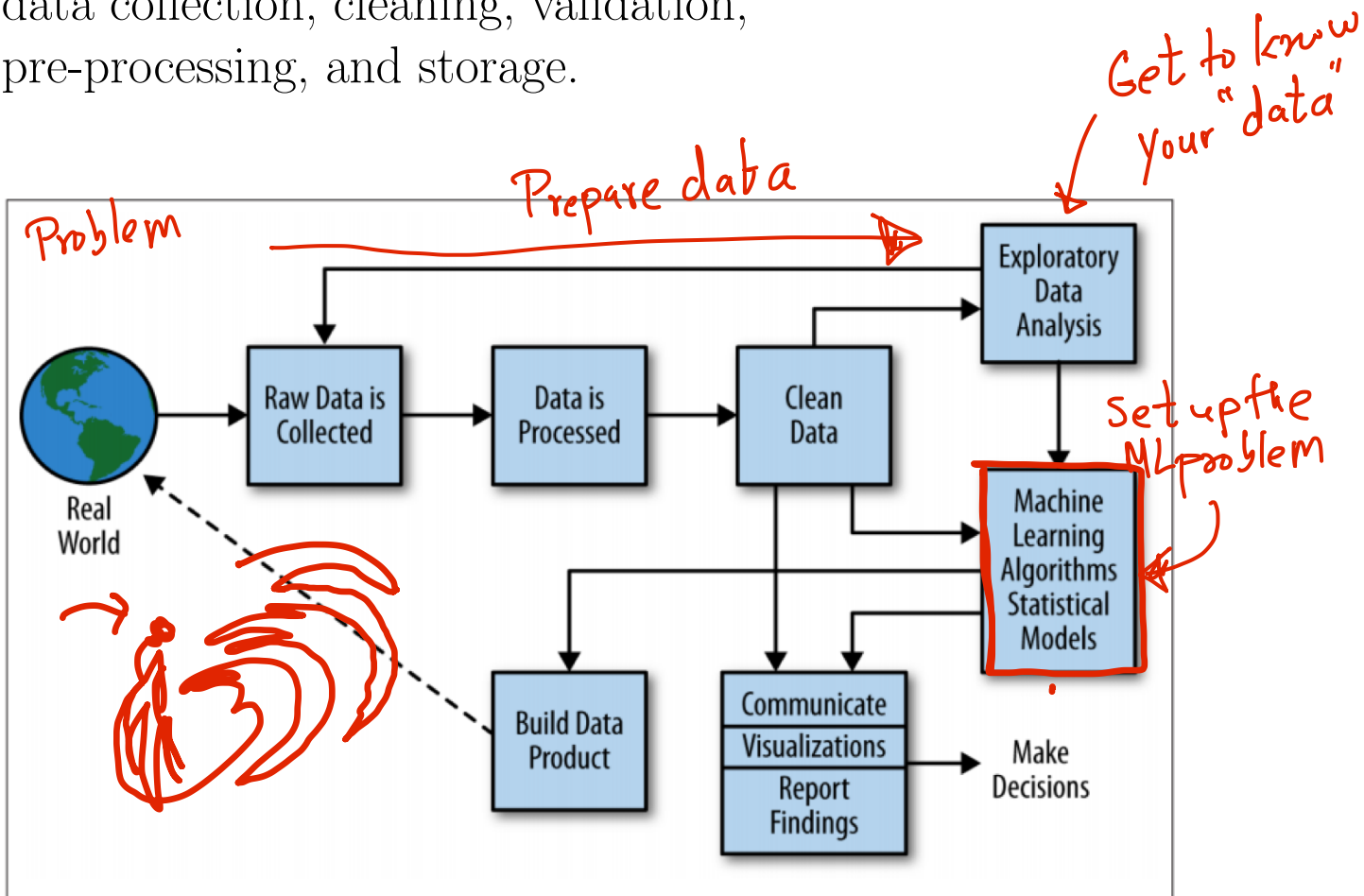
---

<sup>1</sup>Some figures are taken from Hastie, Tibshirani, and Friedman's book on statistical learning and also from Chris Bishop's Machine learning book

# 1 Before Machine Learning

## Acquiring Data

Data is the most important component of modern Machine Learning. There are many important steps that can have a huge impact on the performance of a machine-learning system. To name a few: data collection, cleaning, validation, pre-processing, and storage.



Picture taken from "Doing data science".

# Defining an ML problem

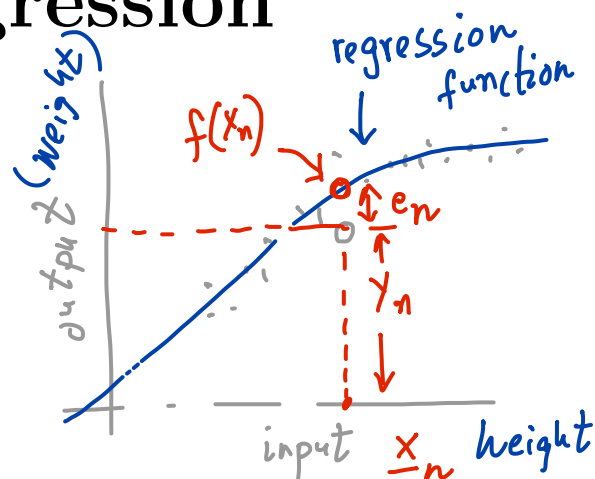
Once we have some data, the next step is to re-define the real-world problem in the context of data, and then to convert it to a machine-learning problem.

ML problems can be categorized into 3 main types: <sup>①</sup> supervised, <sup>②</sup> unsupervised, and <sup>③</sup> reinforcement learning. In practice, a successful end-to-end system might require a combination of these problems.

# 2 ML Problem: Regression

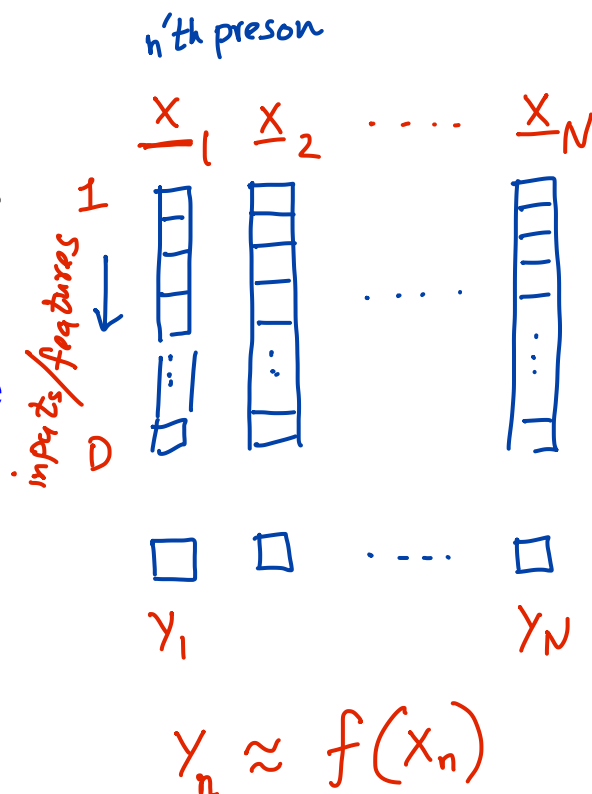
## What is regression?

Regression is to relate input variables to the output variable, to either predict outputs for new inputs and/or to understand the effect of the input on the output.

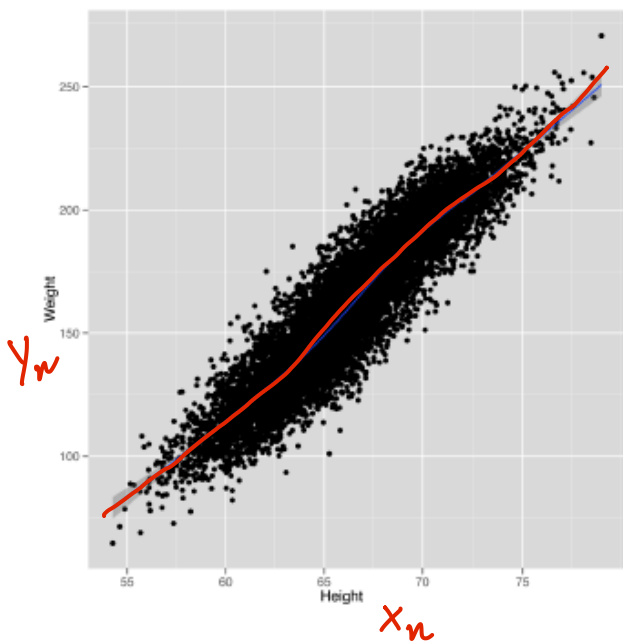


## Dataset for regression

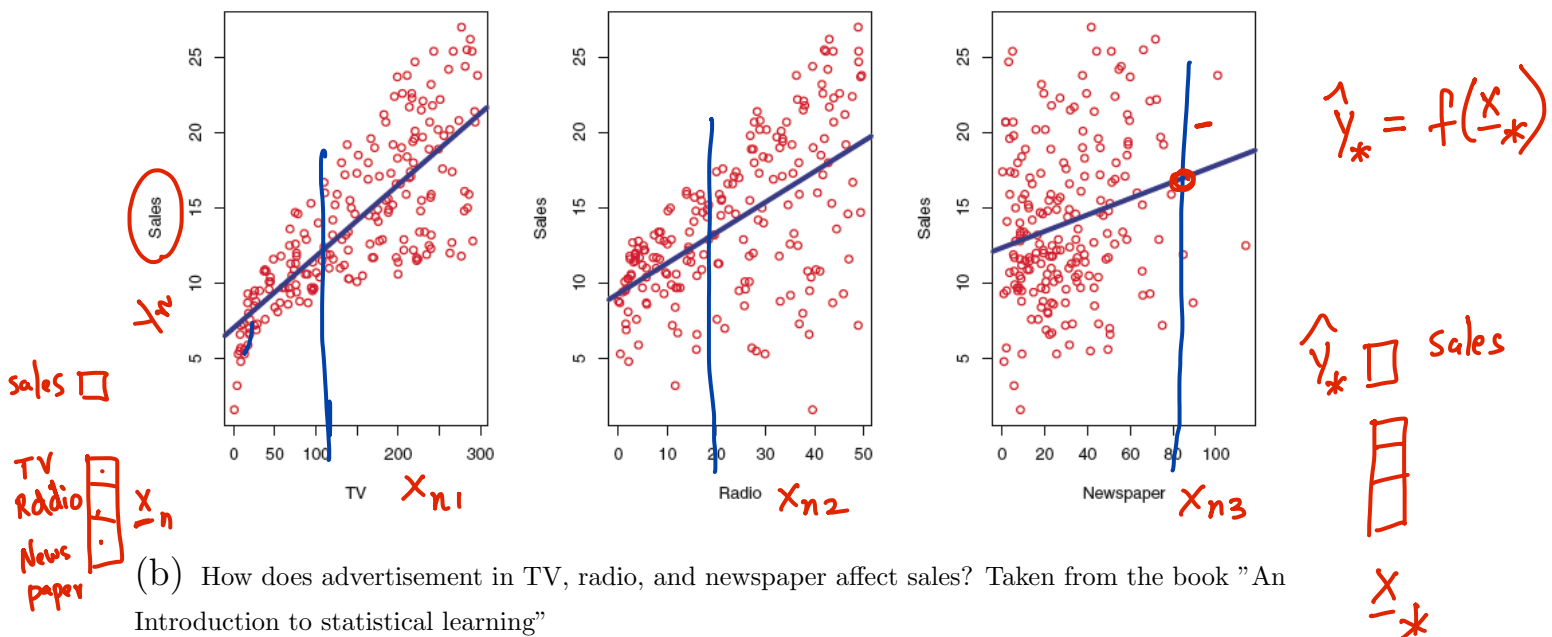
In regression, data consists of pairs  $(y_n, \mathbf{x}_n)$ , where  $y_n$  is the  $n$ 'th output and  $\mathbf{x}_n$  is a vector of  $D$  inputs. Number of pairs  $N$  is the data-size and  $D$  is the dimensionality.



## Examples of regression



(a) Height is correlated with weight. Taken from “Machine Learning for Hackers”



## Two goals of regression

In **prediction**, we wish to predict the output for a new input vector, e.g. what is the weight of a person who is 170 cm tall?

In **interpretation**, we wish to understand the effect of inputs on output, e.g. are taller people heavier too?

## The regression function


For both the goals, we need to find a function that approximates the output “well enough” given inputs.

$$y_n \approx f(\mathbf{x}_n), \text{ for all } n$$

# Additional Notes

## Prediction vs Interpretation

Some questions to think about: are these prediction tasks or interpretation task?

- Prediction* 1. What is the life-expectancy of a person who has been smoking for 10 years?
- interpretation* 2. Does smoking cause cancer? 
- Prediction*  
*interp...* 3. When the number of packs a smoker smokes per day doubles, their predicted life span gets cut in half? ??
- I, P?* 4. A massive<sup>?</sup> scale earthquake will occur in California within next 30 years.
5. More than 300 bird species in north America could reduce their habitat by half or more by 2080.

[Qualitative vs Quantitative  
Predictive models > Interpretive models

Nate Silver "the<sup>?</sup> Signal & noise ... " ??

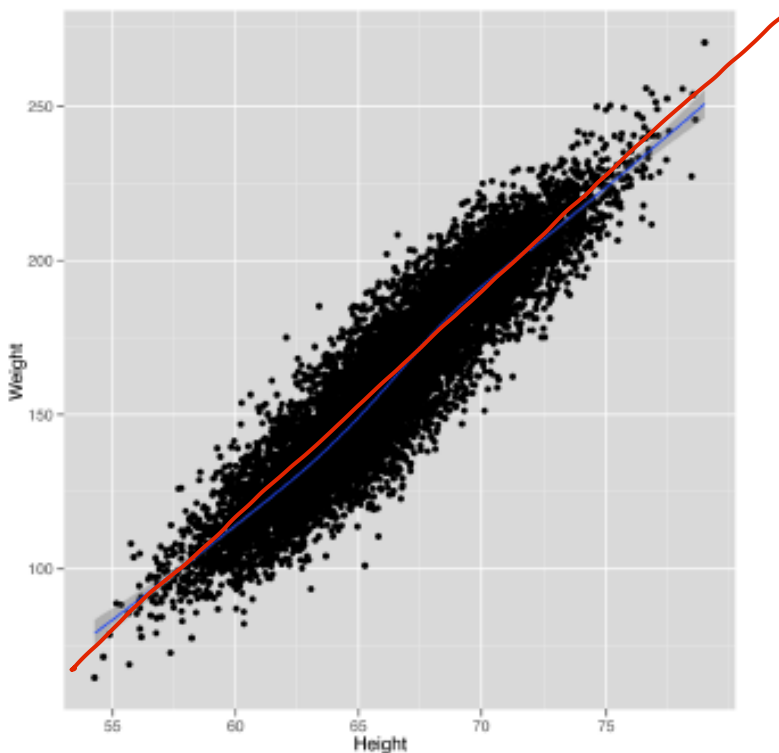
# 3 Model: Linear Regression

## What is it?

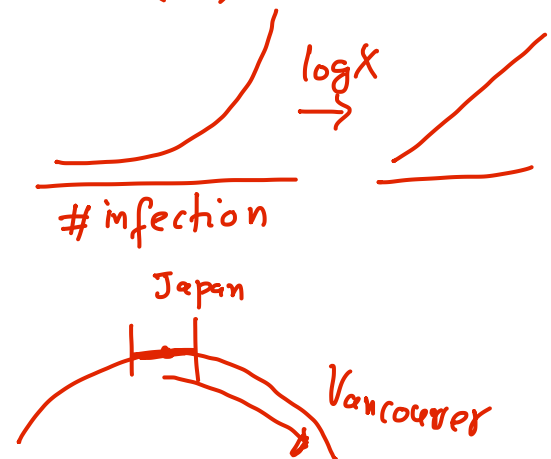
Linear regression is a model that assumes a linear relationship between inputs and the output.

$$y_n \approx f(x_n)$$

↓  
regression function



$ax$  is linear in  $x$   
 $ax^2$  is quadratic in  $x$   
but linear in  $x^2$   
 $a(x')$  with  $x' = x^2$



## Why learn about *linear* regression?

Plenty of reasons: simple, easy to understand, most widely used, easily generalized to non-linear models. Most importantly, you can learn almost all fundamental concepts of ML with regression alone.



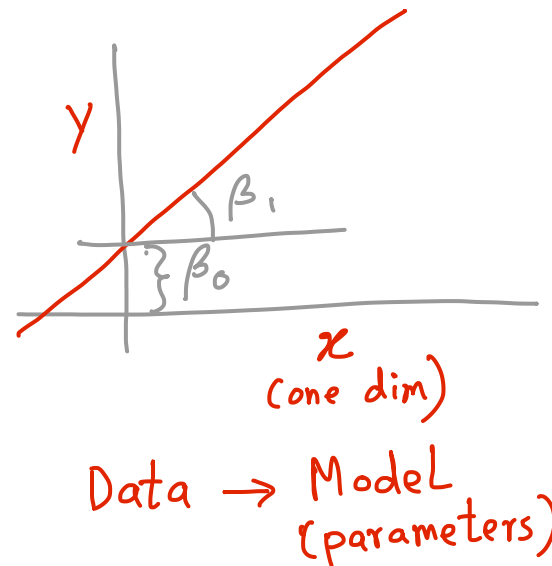
# Simple linear regression

With only one input dimension, it is simple linear regression.

$$y_n \approx f(\mathbf{x}_n) := \beta_0 + \beta_1 x_{n1}$$

"bias"      slope

Here,  $\beta_0$  and  $\beta_1$  are parameters of the model.

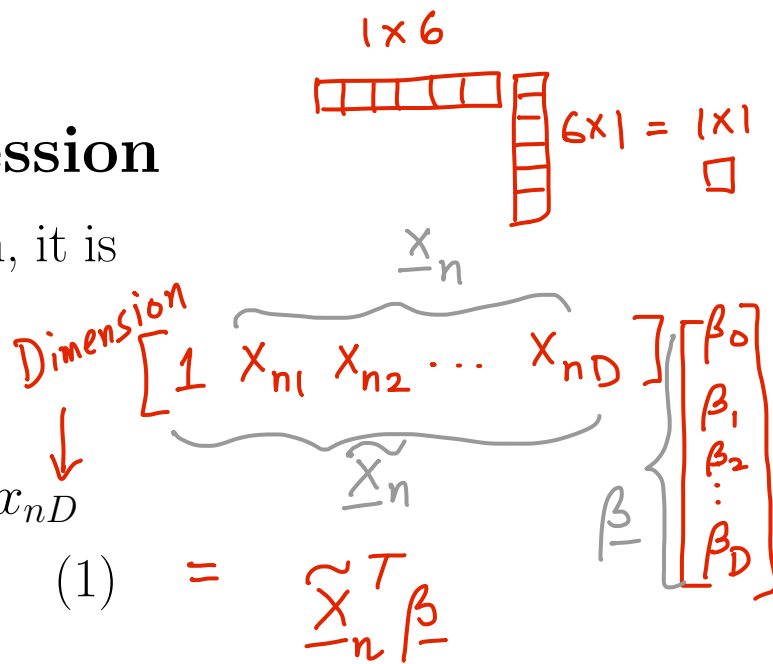


# Multiple linear regression

With multiple input dimension, it is multiple linear regression.

$$y_n \approx f(\mathbf{x}_n) \begin{matrix} \text{bias} \\ \swarrow \end{matrix} := \beta_0 + \beta_1 x_{n1} + \dots + \beta_D x_{nD}$$

$$= \tilde{\mathbf{x}}_n^T \boldsymbol{\beta} \quad \uparrow \quad \uparrow \quad (1)$$



# Learning/estimation/fitting

Given data, we would like to find  $\boldsymbol{\beta} = [\beta_0, \beta_1, \dots, \beta_D]$ . This is called **learning** or **estimating** the parameters or **fitting** the model.

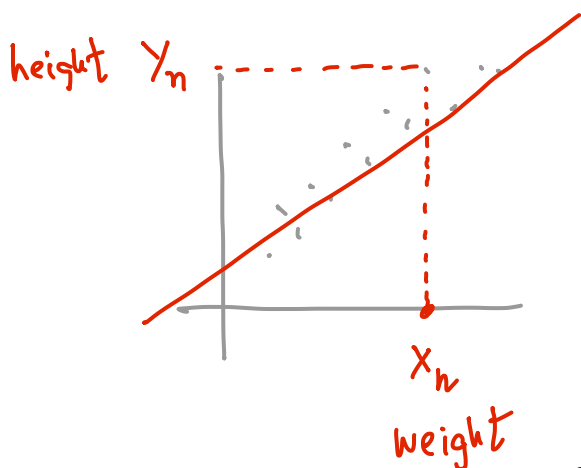
$$A B = C$$

$$N \times D \quad D \times M \quad N \times M$$

$$\longleftrightarrow$$

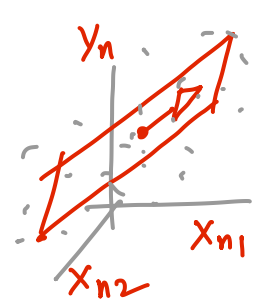
$$1 \times 6 \quad 6 \times 1 \quad 1 \times 1$$

$$y_n \quad \beta_0 \quad \beta_1 \quad x_{n,1} \quad x_{n,2}$$



$$y_1 \quad 180 \approx 0 + 3 \times 60 + 1 \times 2k$$

$$y_2 \quad 185 \approx 0 + 3 \times 70 + 1 \times 3k$$



# Additional Notes

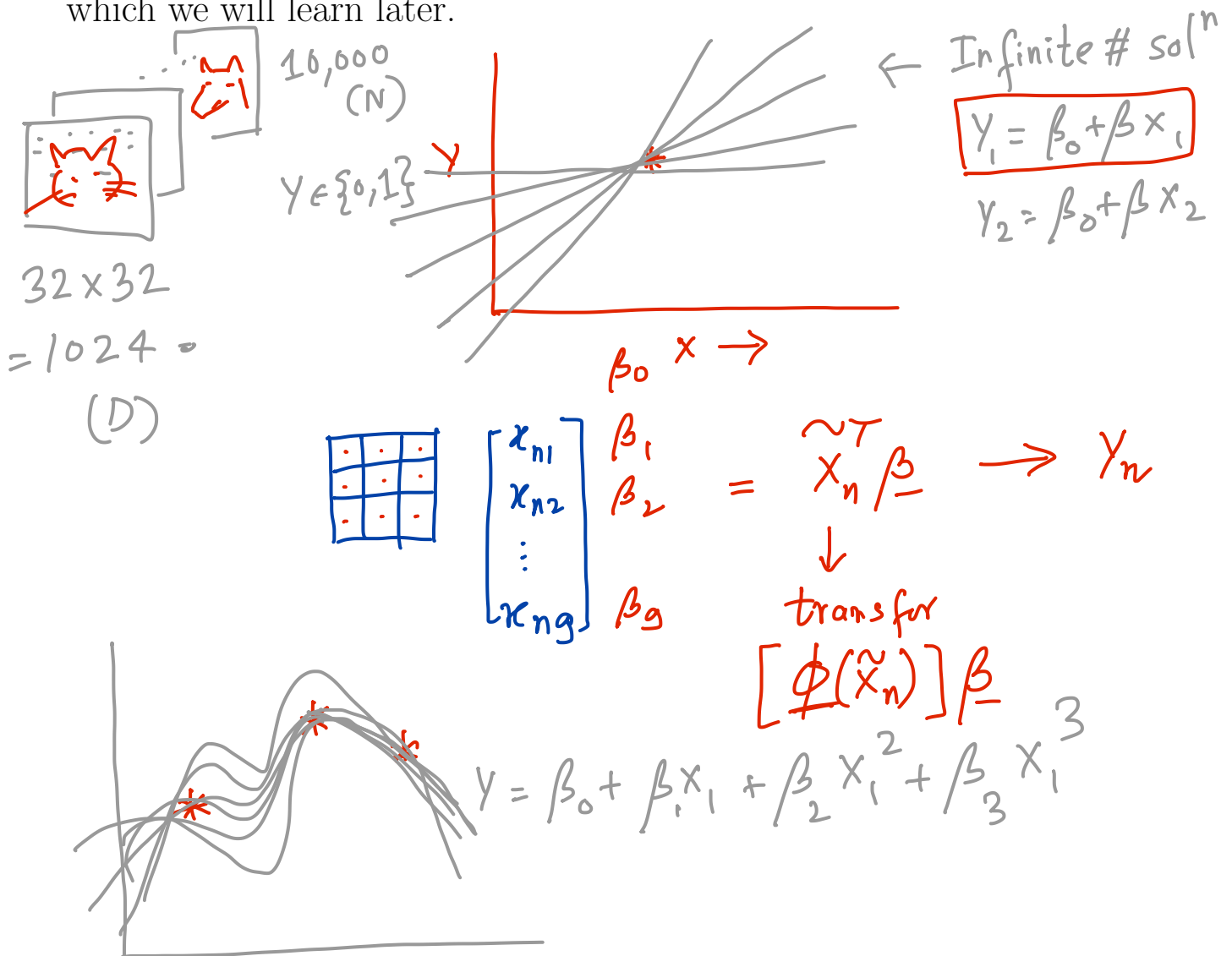
## $p > n$ Problem

Consider the following simple situation: You have  $N = 1$  and you want to fit  $y_1 \approx \beta_0 + \beta_1 x_{11}$ , i.e. you want to find  $\beta_0$  and  $\beta_1$  given one pair  $(y_1, x_{11})$ . Is it possible to find such a line? *Any line is good!*

*degrees of freedom,  $(N, D)$*

This problem is related to something called  $p > n$  problem. In our notation, this will be called  $D > N$  problem, i.e. the number of parameters exceeds number of data examples.

Similar issues will arise when we use gradient descent or least-squares to fit a linear model. These problems are all solved by using regularization, which we will learn later.



# 4 Cost Function: MSE

## Motivation

Consider the following models.

1-parameter model:  $y_n \approx \beta_0$

2-parameter model:  $y_n \approx \beta_0 + \beta_1 x_{n1}$

How can we estimate (or guess) values of  $\beta$  given the data  $D$ ?

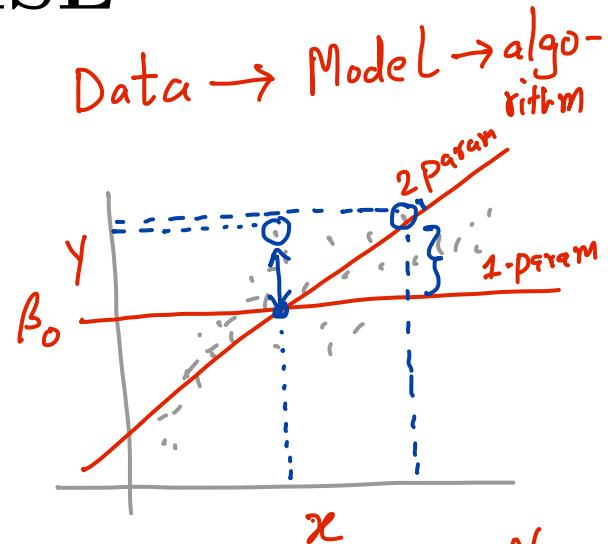
## What is a cost function?

Cost functions (or utilities or energy) are used to learn parameters that explain the data well. They define how costly our mistakes are.

## Two desirable properties of cost functions

When  $y$  is real-valued, it is desirable that the cost is symmetric around 0, since both +ve and -ve errors should be penalized equally.

Also, our cost function should penalize “large” mistakes and “very-large” mistakes almost equally.



$$\text{Data: } D = \{y_n, x_n\}_{n=1}^N$$

$$\text{Model: } y_n \approx \tilde{x}_n^T \beta$$

$$\text{Cost func: } C(D, \beta)$$

$$\text{Error: } e_n = y_n - \tilde{x}_n^T \beta$$



$$C(e_1, e_2, \dots, e_N)$$

# Mean Square Error (MSE)

MSE is one of the most popular cost function.

$$C(e_1, e_2, \dots, e_N) = \frac{1}{2N} \sum_{i=1}^N \text{MSE}(e_i) = \frac{1}{2N} \sum_{i=1}^N e_i^2$$

$$\text{MSE}(\beta) := \frac{1}{2N} \sum_{n=1}^N [y_n - f(\mathbf{x}_n)]^2$$

Assumption: "all data are equal"  
More precisely: i.i.d. (independent & identically distributed)

Does it have both the properties?  
symmetric, yes!

## An exercise for MSE

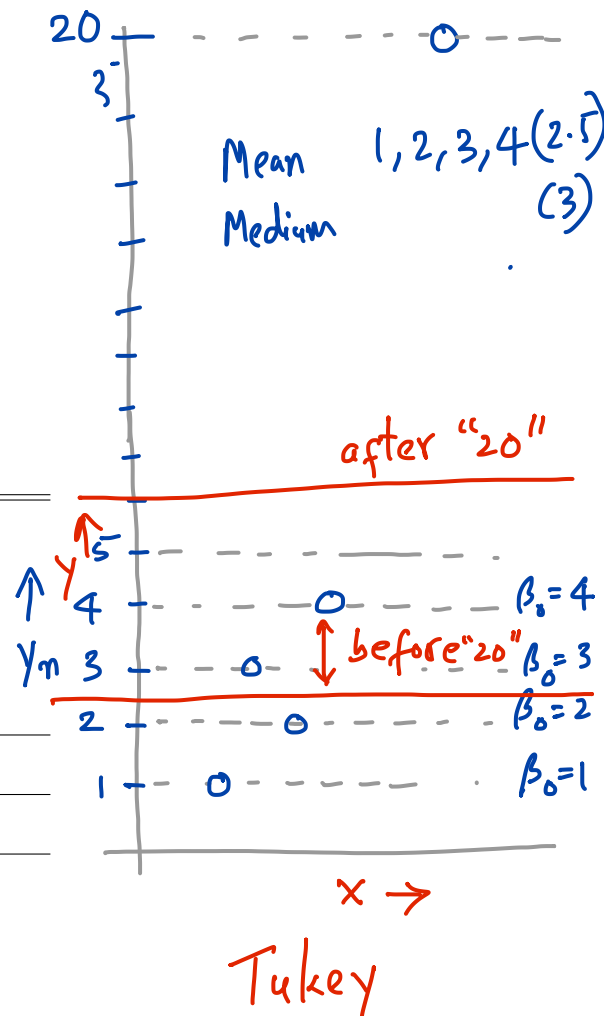
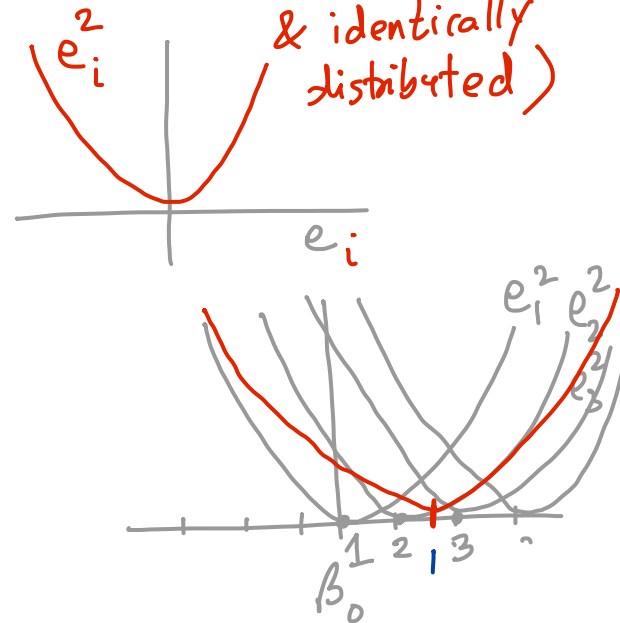
Compute MSE for 1-param model:

$$\underline{\mathcal{L}(\beta_0)} := \frac{1}{2N} \sum_{n=1}^N [y_n - \beta_0]^2 \quad (2)$$

Each row contains a  $y_n$  and column is  $\beta_0$ . First, compute MSE for for  $y_n = \{1, 2, 3, 4\}$  and draw MSE as a function of  $\beta_0$  (by adding the first four rows). Then add  $y_n = 20$  to it, and redraw MSE. What do you observe and why?

$y_n \backslash \beta_0$	1	2	3	4	5	6	7
$e_1$	0	1	2	3	4		
$e_2$	1	0	1	2	3		
$e_3$	2	1	0	1	2		
$e_4$	3	2	1	0	1		
MSE	14	6	6	14	30		
20	19	18	17	16	15		
MSE	.	.	.	.	.		

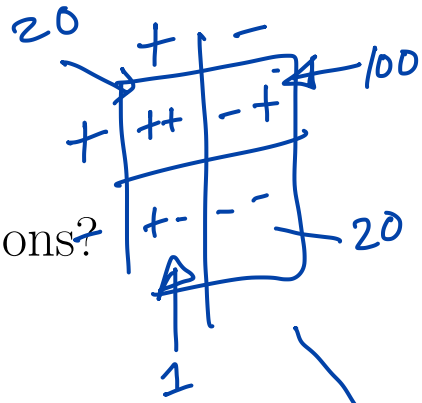
it may not be "bad", scale of data



# Additional Notes

## A question for cost functions

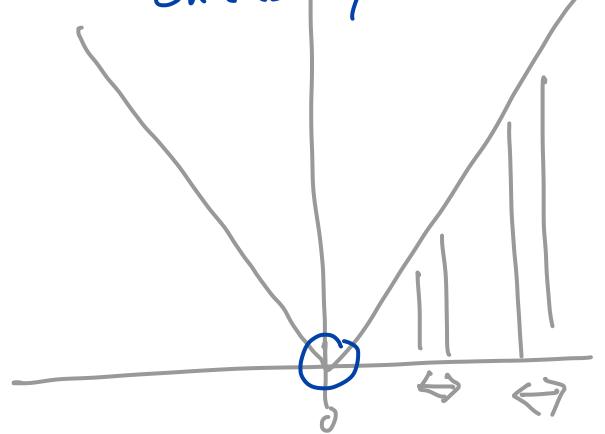
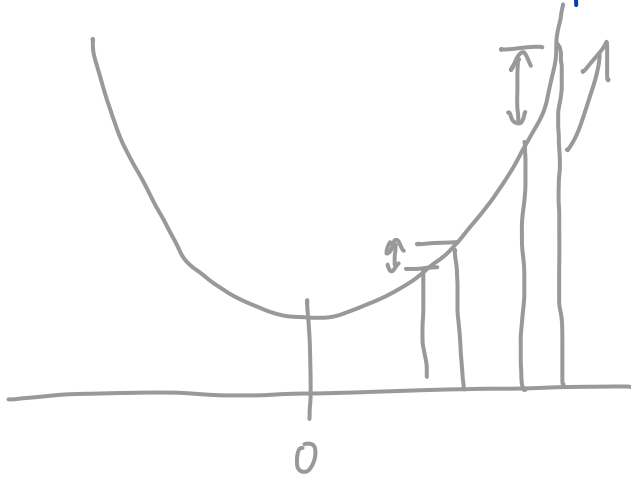
Is there an automatic way to define loss functions?



## Nasty cost functions: Visualization

See Andrej Karpathy Tumblr post for many cost functions gone "wrong" for neural network. <http://lossfunctions.tumblr.com/>.

"Cost-sensitive" Cost functions: *loss/* Need to design them on case-by-case basis.

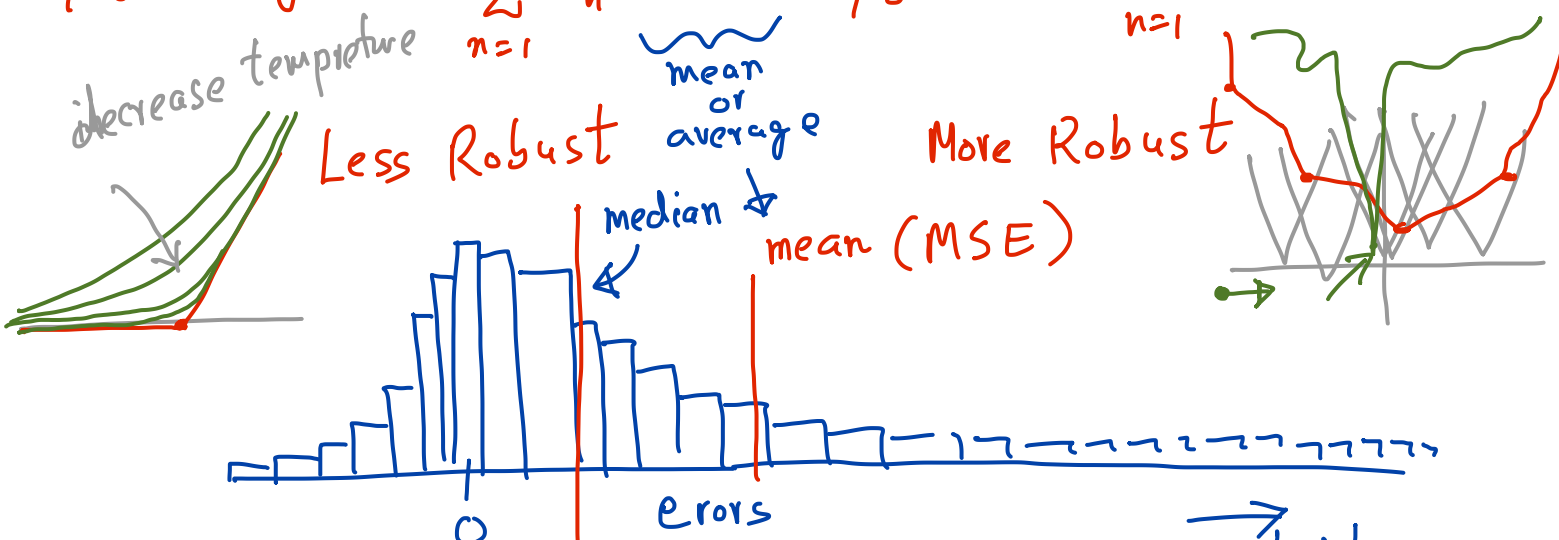


$$\beta_0 = \arg \min \sum_{n=1}^N e_n^2 = \text{mean} \left( \begin{matrix} y_1 \\ \vdots \\ y_N \end{matrix} \right) \quad \beta_0 = \arg \min \sum_{n=1}^N |e_n| = \text{median} \left( \begin{matrix} y_1 \\ \vdots \\ y_N \end{matrix} \right)$$

decrease temperature

Less Robust

More Robust



Desirable

Properties of cost function  $|e|$ : Continuous, differentiable (twice), Convex (strictly, strongly)

# 5 Algorithm 1: Gradient Descent

## Learning/estimation/fitting

Given a cost function  $\mathcal{L}(\beta)$ , we wish to find  $\beta^*$  that minimizes the cost:

$$\rightarrow \min_{\beta} \mathcal{L}(\beta), \quad \text{subject to } \underline{\beta} \in \mathbb{R}^{D+1}$$

This is learning posed as an optimization problem. We will use an algorithm to solve the problem.

## Grid search

$$\begin{bmatrix} | & | & | & | \\ \hline & & & \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix} = \tilde{X}_n^T \underline{\beta}$$

Grid search is one of the simplest algorithms where we compute cost over a grid (of say  $M$  points) to find the minimum. This is extremely simple and works for any kind of loss when we have very few parameters and the loss is easy to compute.

For a large number of parameters, however, grid search has too many "for-loops", resulting in exponential computational complexity. Choosing a good range of values is another problem.

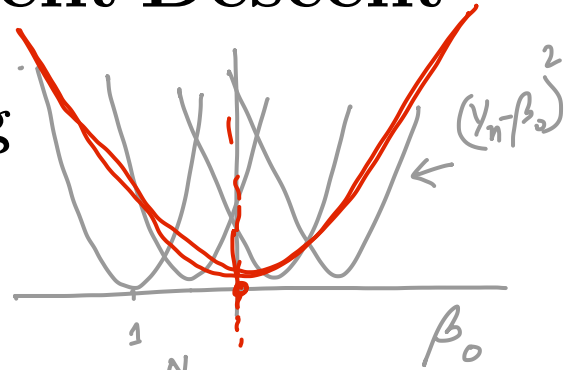
Q: What is the computational complexity of this algorithm?

Are there any other issues?

Look into Wikipedia (big-O notation)

$$O(ND M^D) \leftarrow \text{All the points } \underline{\beta}$$

13 for each MSE[i]

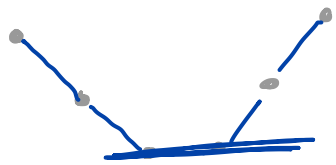


$$\mathcal{L}(\beta_0) = \frac{1}{2N} \sum_{n=1}^N (y_n - \beta_0)^2$$

$$\mathcal{L}(\beta_0, \beta_1) = \frac{1}{2N} \sum_{n=1}^N (y_n - \beta_0 - \beta_1 x_{n1})^2$$

$$\mathcal{L}\left(\begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}\right) = \frac{1}{2N} \sum_{n=1}^N (y_n - \tilde{X}_n^T \underline{\beta})^2$$

$$\mathcal{L}(\underline{\beta}) : \mathbb{R}^{D+1} \rightarrow \mathbb{R}$$



for each dimension 1:D  
say M points

$$\beta_0 \leftarrow [-10; 0.1; 10]$$

for i = 1: length(beta\_0)

MSE[i] ← compute loss(beta\_0, D)

end

i ← argmin(MSE)

MSE ← 0

for n = 1: N

$e_n \leftarrow y_n - \tilde{X}_n^T \underline{\beta}$

MSE = MSE +  $e_n^2$

end

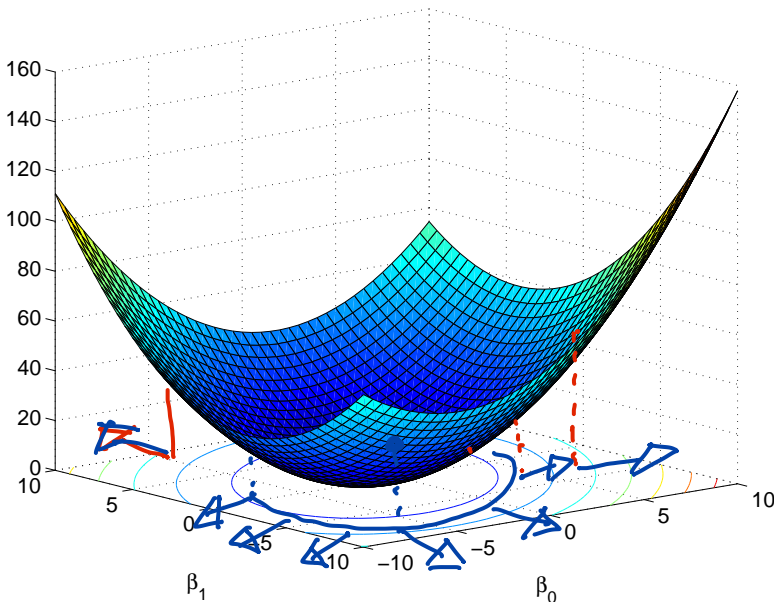
$O(ND)$  ~~MD~~ Loss function is "smooth". (Lipschitz)

## Follow the gradient

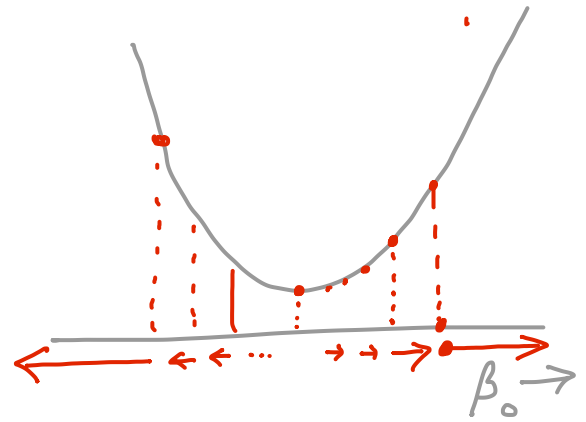
A gradient (at a point) is the slope of the tangent (at that point). It points to the direction of largest increase of the function.

For 2-parameter model, MSE is shown below.

(I used  $\mathbf{y}^T = [2, -1, 1.5]$  and  $\mathbf{x}^T = [-1, 1, -1]$ ).



$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta}$$



$$\begin{aligned} \frac{\partial \mathcal{L}(\beta_0)}{\partial \beta_0} &= \frac{1}{2N} \frac{\partial}{\partial \beta_0} \sum_{n=1}^N (y_n - \beta_0)^2 \\ &= \frac{1}{2N} \sum_{n=1}^N -2(y_n - \beta_0) \\ &= \frac{1}{N} \sum_{n=1}^N (\beta_0 - y_n) \end{aligned}$$

error =  $e_n$

$$\mathbb{R}^2 \ni \frac{\partial \mathcal{L}(\beta_0, \beta_1)}{\partial (\beta_0, \beta_1)} = \frac{\partial}{\partial \underline{\beta}} \frac{1}{2N} \sum_{n=1}^N (y_n - \beta_0 - \beta_1 x_{n1})^2$$

error =  $e_n$



$$\begin{aligned} &= \frac{1}{2N} \sum_{n=1}^N \begin{bmatrix} -2(y_n - \beta_0 - \beta_1 x_{n1}) \\ -2(y_n - \beta_0 - \beta_1 x_{n1}) x_{n1} \end{bmatrix} = \begin{bmatrix} \frac{1}{N} \sum_{n=1}^N (\tilde{\mathbf{x}}_n^T \underline{\beta} - y_n) \\ \frac{1}{N} \sum_{n=1}^N (\tilde{\mathbf{x}}_n^T \underline{\beta} - y_n) x_{n1} \end{bmatrix} \\ &= \frac{1}{N} \sum_{n=1}^N e_n \begin{bmatrix} 1 \\ x_{n1} \end{bmatrix} = \frac{1}{N} \sum_{n=1}^N e_n \tilde{\mathbf{x}}_n \end{aligned}$$

error (scalar)

error      input

# Batch gradient descent

To minimize the function, take a step in the (opposite) direction of the gradient

$$\underbrace{\beta^{(k+1)}}_{\text{New}} \leftarrow \underbrace{\beta^{(k)}}_{\text{current}} - \alpha \frac{\partial \mathcal{L}(\beta^{(k)})}{\partial \beta}$$

where  $\alpha > 0$  is the step-size (or learning rate).

We want the sequence  $\beta^{(0)} \beta^{(1)} \beta^{(2)} \dots$  to "converge" to  $\beta^{(*)}$

Gradient descent for 1-parameter model to minimize MSE:

$$\beta_0^{(k+1)} = (1 - \alpha)\beta_0^{(k)} + \alpha \bar{y}$$

Where  $\bar{y} = \sum_n y_n / N$ . <sup>Q:</sup> When is this sequence guaranteed to converge?

(for what value of  $\alpha$ ?)  $\leftarrow$

## Gradients for MSE

$$\mathcal{L}(\beta) = \frac{1}{2N} \sum_{n=1}^N (y_n - \tilde{\mathbf{x}}_n^T \beta)^2 \quad (3)$$

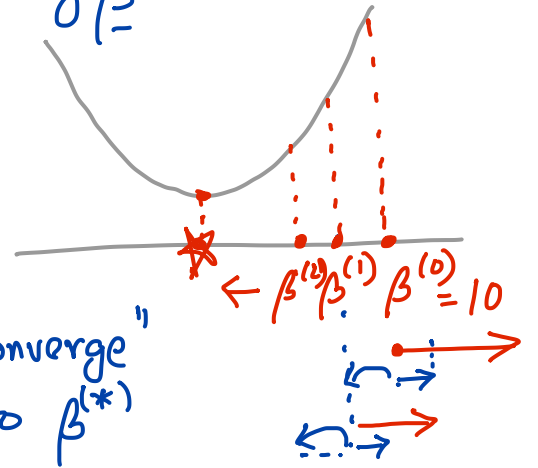
then the gradient is given by,

$$\frac{\partial \mathcal{L}}{\partial \beta} = -\frac{1}{N} \sum_{n=1}^N \underbrace{(y_n - \tilde{\mathbf{x}}_n^T \beta)}_{\text{error} = e_n} \tilde{\mathbf{x}}_n \quad (4)$$

What is the computational complexity of batch gradient descent?

$$\frac{\partial \mathcal{L}(\beta)}{\partial \beta} \Big|_{\beta = \beta^{(k)}} \leftarrow k\text{'th iteration}$$

$$= \frac{\partial \mathcal{L}(\beta^{(k)})}{\partial \beta}$$



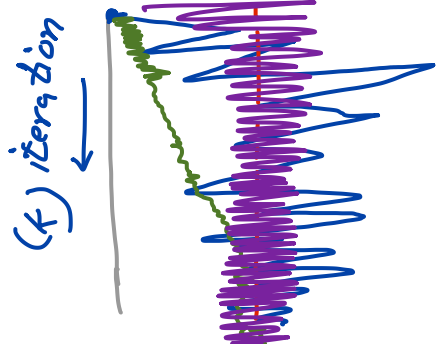
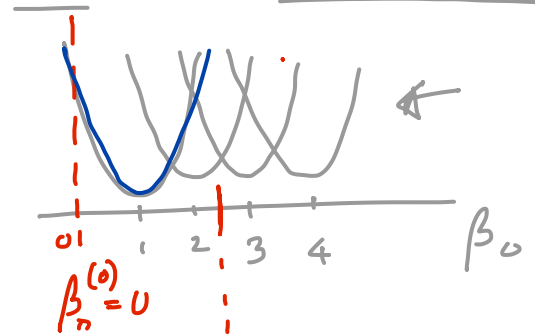
$$\frac{\partial \mathcal{L}(\beta_0)}{\partial \beta_0} = \frac{1}{N} \sum_{n=1}^N (\beta_0 - y_n)$$

$$= \beta_0 - \frac{1}{N} \sum_{n=1}^N y_n$$

$$= \beta_0 - \bar{y}$$

$$\beta_0^{(k+1)} \leftarrow \beta_0^{(k)} - \alpha (\beta_0^{(k)} - \bar{y})$$

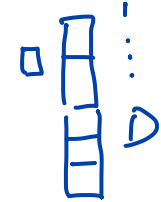
$$= (1 - \alpha)\beta_0^{(k)} + \alpha \bar{y}$$





$$O(N D M^D) \xrightarrow{\text{grid search}} O(N D) \xrightarrow{\text{GD}} O(D)$$

# Stochastic gradient descent

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \beta} &= \frac{1}{N} \sum_{n=1}^N e_n \underline{x}_n \\ &= \frac{1}{N} N \times e_i \underline{x}_i \\ &= e_i \underline{x}_i \end{aligned}$$


When  $N$  is large, choose a random pair  $(\mathbf{x}_i, y_i)$  in the training set and approximate the gradient:

$$\frac{\partial \mathcal{L}}{\partial \beta} \approx -\frac{1}{N} \left[ N(y_i - \tilde{\mathbf{x}}_i^T \beta) \tilde{\mathbf{x}}_i \right] \quad (5)$$

Using the above “stochastic” gradient, take a step:

$$\beta^{(k+1)} = \beta^{(k)} + \alpha^{(k)} (y_i - \tilde{\mathbf{x}}_i^T \beta^{(k)}) \tilde{\mathbf{x}}_i$$

What is the computational complexity?

For convergence,  $\alpha^k \rightarrow 0$  “appropriately”. One such condition called Robbins-Monroe condition suggests to take  $\alpha^k$  such that:

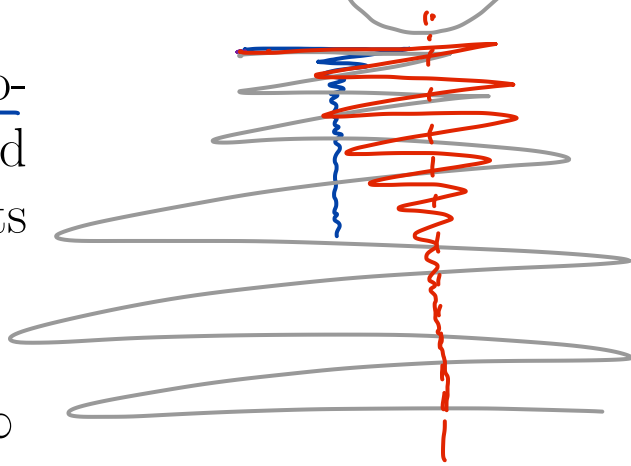
$$\sum_{k=1}^{\infty} \alpha^{(k)} = \infty, \quad \sum_{k=1}^{\infty} (\alpha^{(k)})^2 < \infty \quad (6)$$

One way to obtain such sequence is  $\alpha^{(k)} = \frac{1}{(1+k)^r}$  where  $r \in (0.5, 1)$ .

$$\left\{ \begin{array}{l} r=2 \\ r=0.4 \end{array} \right\}$$

Assumption: i.i.d. over  $n$

For constant learning rate, SGD does “not” converge.



$$\alpha^{(k)} < \text{constant}$$

$$\sum \alpha^{(k)} = \infty$$

$$\sum (\alpha^{(k)})^2 < \infty$$

Data  
 $\{ \underline{x}_i, y_i \}_{i=1}^N$

Model  
 $f(\underline{y}_n) = \underline{\tilde{x}}_n^T \underline{\beta}$

Algorithm  
 $\min_{\underline{\beta}} \frac{1}{2N} \sum_{n=1}^N (y_n - \underline{\tilde{x}}_n^T \underline{\beta})^2$

$$\frac{\partial \Delta(\underline{\beta}^{(k)})}{\partial \underline{\beta}} = \frac{1}{N} \sum_{n=1}^N (y_n - \underbrace{f_{\underline{\beta}^{(k)}}(\underline{x}_i)}_{\text{forward pass}}) \underline{x}_i$$

$$\begin{bmatrix} 1 \\ \underline{x}_n \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix}$$

$$\underline{\tilde{x}}_n^T \underline{\beta}^{(k)}$$

↓

$$e_n$$

↓

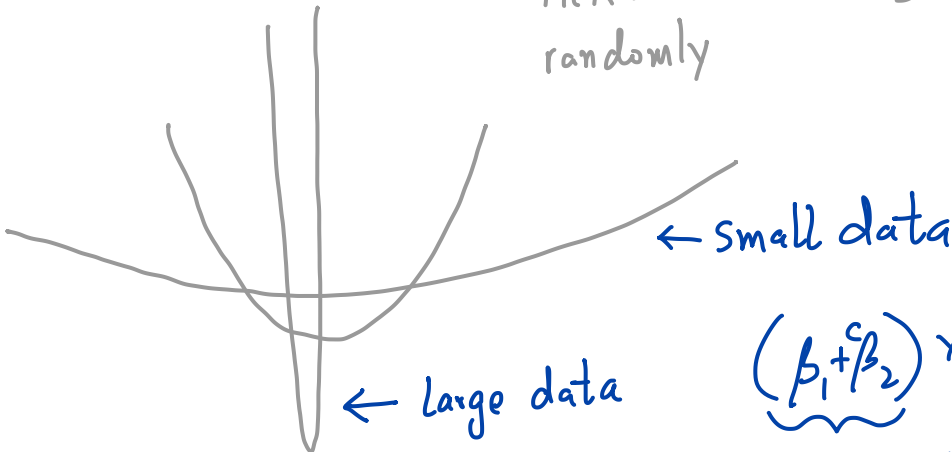
$$e_n \underline{x}_i$$

$$\approx e_i \underline{x}_i$$

↑  
Pick  $i$   
randomly

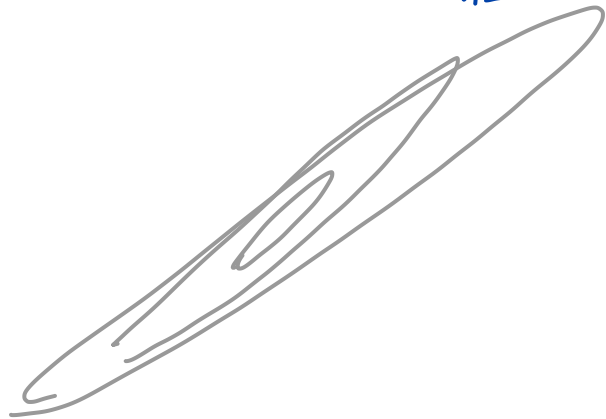
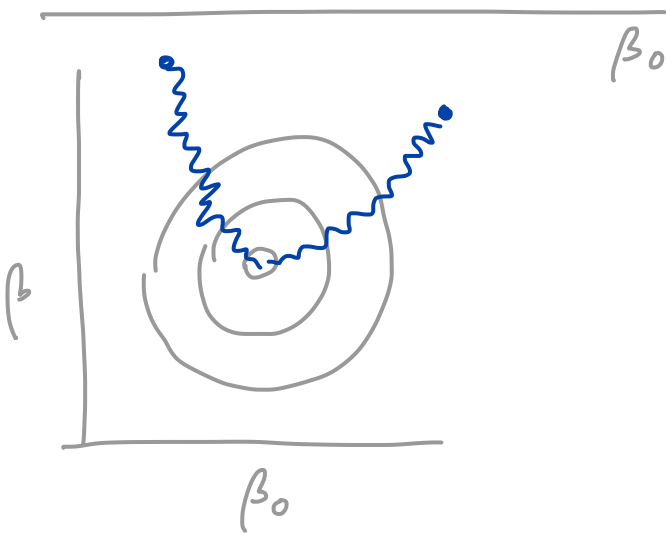
compute error

back pass



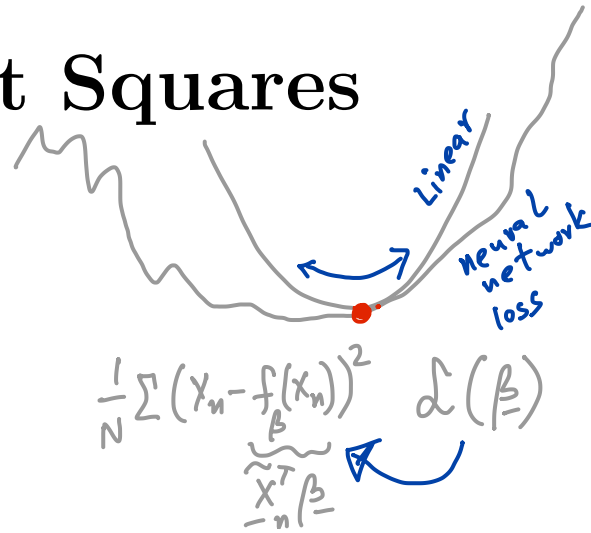
$$\underbrace{(\beta_1 + c\beta_2)}_{\beta_{\text{new}} \underline{x}_{n1}} \underline{x}_{n1} = \beta_1 \underline{x}_{n1} + \beta_2 \underline{x}_{n2}$$

$$\underline{x}_{n2} = c \underline{x}_{n1}$$



# 6 Algorithm 2: Least Squares

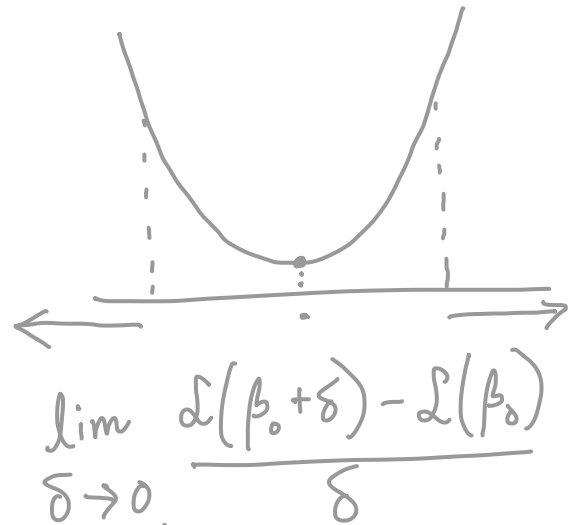
In rare cases, we can compute the minimum of the cost function analytically. Linear regression using MSE is one such case. The solution is obtained using normal equations. This is called least squares. (Gauss)



To derive the equation, we use the stationary optimality conditions. See the lecture notes for Gradient Descent.

$$\frac{\partial \mathcal{L}(\beta^*)}{\partial \beta} = 0$$

$\beta^{(0)} \quad \beta^{(1)} \quad \dots \quad \beta^*$   
 $\xrightarrow{\text{converges}}$



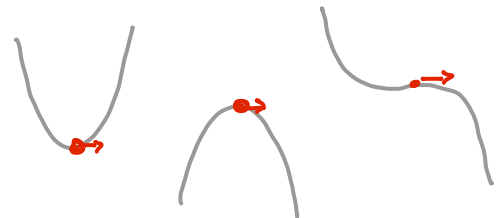
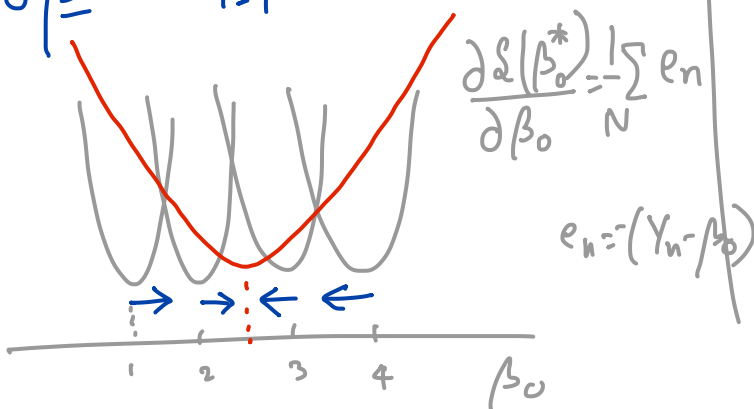
Using this, derive the normal equation for 1-parameter model.

$$\frac{\partial \mathcal{L}(\beta^*)}{\partial \beta} = \frac{1}{2N} \sum_{i=1}^N e_n \tilde{x}_n = 0$$

Informal statement  
 When  $\beta = \underset{\beta}{\operatorname{argmin}} \mathcal{L}(\beta)$

then  $\frac{\partial \mathcal{L}(\beta^*)}{\partial \beta} = 0$

(Opposite is not true)



$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha (y_i - \beta_0^{(k)})$$

# Normal equations

Recall the expression of the gradient for multiple linear regression:

$$\frac{\partial \mathcal{L}}{\partial \beta} = -\frac{1}{N} \tilde{\mathbf{X}}^T \mathbf{e} = -\frac{1}{N} \tilde{\mathbf{X}}^T (\mathbf{y} - \tilde{\mathbf{X}}\beta)$$

Set it to zero to get the normal equations for linear regression.

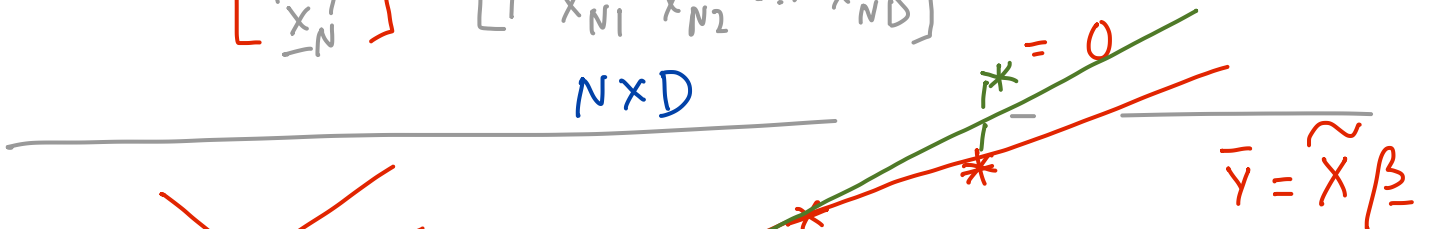
$$\tilde{\mathbf{X}}^T \mathbf{e} = \tilde{\mathbf{X}}^T (\mathbf{y} - \tilde{\mathbf{X}}\beta) = 0$$

implying that the error is orthogonal to rows of  $\tilde{\mathbf{X}}^T$  and columns of  $\tilde{\mathbf{X}}$ .

$$\tilde{\mathbf{X}} = \begin{bmatrix} \tilde{x}_1^T \\ \vdots \\ \tilde{x}_N^T \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1D} \\ 1 & x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix}$$

$N \times D$

$$\begin{aligned} \frac{\partial \mathcal{L}(\beta^*)}{\partial \beta} &= \frac{1}{N} \sum_{n=1}^N e_n \tilde{x}_n \\ &= \frac{1}{N} \left( e_1 \begin{bmatrix} 1 \\ \tilde{x}_1 \end{bmatrix} + \dots + e_N \begin{bmatrix} 1 \\ \tilde{x}_N \end{bmatrix} \right) \\ &= \frac{1}{N} \left[ \begin{bmatrix} 1 \\ \tilde{x}_1 \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} \dots \begin{bmatrix} 1 \\ \tilde{x}_N \end{bmatrix} \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} \right] \\ &\quad \underbrace{\hspace{10em}}_{(D+1) \times N} \quad \underbrace{\hspace{5em}}_{N \times 1} \\ &= \frac{1}{N} \tilde{\mathbf{X}}^T \mathbf{e} \\ &= \frac{1}{N} \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}}\beta^* - \mathbf{y}) \end{aligned}$$



~~$$e_1 \begin{bmatrix} 1 \\ x_{11} \end{bmatrix} = 0$$~~

$$e_1 \begin{bmatrix} 1 \\ x_{11} \end{bmatrix} + e_2 \begin{bmatrix} 1 \\ x_{21} \end{bmatrix} = 0$$

$$+ e_3 \begin{bmatrix} 1 \\ x_{31} \end{bmatrix} = 0$$

$$\underbrace{(\beta_0 + \beta_1 x_{11} - y_1)}_{e_1 = 0} \begin{bmatrix} 1 \\ x_{11} \end{bmatrix} = 0$$

$\tilde{\mathbf{X}} = \begin{bmatrix} 1 & x_{11} \\ 1 & x_{21} \\ 1 & x_{31} \end{bmatrix} \begin{matrix} \tilde{x}_1 \\ \tilde{x}_2 \\ \tilde{x}_3 \end{matrix}$  concept of rank  
 $\downarrow$   
 Maximum "degree of freedom"

$$\sum_{n=1}^N e_n \tilde{\mathbf{x}}_n = 0 \Rightarrow \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}} \underline{\beta}^* - \underline{\mathbf{y}}) \text{ } \left. \vphantom{\sum_{n=1}^N e_n \tilde{\mathbf{x}}_n = 0} \right\} \text{error}$$

# Geometric Interpretation

Linear combination of  $\underline{\mathbf{x}}_d$

$$\tilde{\mathbf{X}}^T \underline{\mathbf{e}} = 0$$

Denote the  $d$ 'th column of  $\tilde{\mathbf{X}}$  by  $\tilde{\mathbf{x}}_d$ .

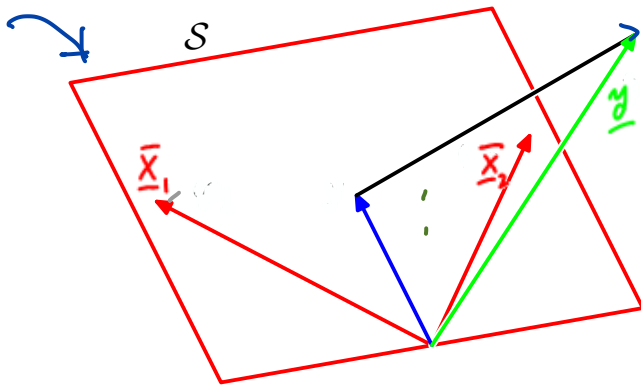
$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \tilde{\mathbf{X}} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1D} \\ 1 & x_{21} & x_{22} & \dots & x_{2D} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{N1} & x_{N2} & \dots & x_{ND} \end{bmatrix}$$

$$\begin{matrix} \tilde{\mathbf{x}}_1^T \\ \tilde{\mathbf{x}}_2^T \\ \vdots \\ \tilde{\mathbf{x}}_N^T \end{matrix}$$

The normal equations suggest to choose a vector in the span of  $\tilde{\mathbf{X}}$ .

The following figure illustrates this (taken from Bishop's book)

space of all vectors obtained with a linear combination of  $\underline{\mathbf{x}}_1, \dots, \underline{\mathbf{x}}_D$



$$\underline{\mathbf{y}} - \tilde{\mathbf{X}} \underline{\beta}$$

$$\underline{\mathbf{e}} = \tilde{\mathbf{X}} \underline{\beta} - \underline{\mathbf{y}}$$

$$\begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} \leftarrow \begin{matrix} \tilde{\mathbf{x}}_1^T \underline{\beta} - y_1 \\ \tilde{\mathbf{x}}_2^T \underline{\beta} - y_2 \\ \vdots \\ \tilde{\mathbf{x}}_N^T \underline{\beta} - y_N \end{matrix}$$

$$\begin{matrix} N \times 1 & N \times D & D \times 1 & N \times 1 \end{matrix}$$

$$\underbrace{\tilde{\mathbf{X}} \quad \underline{\beta}}_{N \times 1}$$

# Least-squares

When  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  is invertible, we have a closed-form expression for the minimum.

$$\underline{\beta}^* = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \underline{\mathbf{y}}$$

$$\tilde{\mathbf{X}}^T (\tilde{\mathbf{X}} \underline{\beta}^* - \underline{\mathbf{y}}) = 0$$

$$\Rightarrow (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}) \underline{\beta}^* = \tilde{\mathbf{X}}^T \underline{\mathbf{y}}$$

We can predict values for a new  $\mathbf{x}_*$ .

$$\hat{y}_* = \underline{\tilde{\mathbf{x}}_*^T} \underline{\beta}^* = \underline{\tilde{\mathbf{x}}_*^T} (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \underline{\mathbf{y}}$$

$$= \underline{\beta}_0 + \beta_1 \underline{x}_1 + \dots + \beta_D \underline{x}_D$$

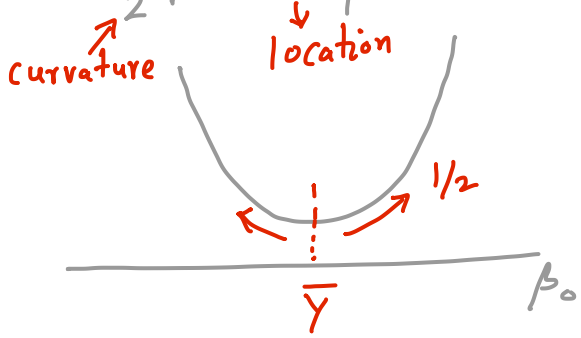
Gram matrix assume that it's invertible

$$\frac{1}{2N} \sum_{i=1}^N (y_i - \beta_0)^2 = \frac{1}{2N} \left( \sum_{i=1}^N \beta_0^2 - 2 \sum_{i=1}^N y_i \beta_0 + \sum_{i=1}^N y_i^2 \right)$$

$$= \frac{1}{2} \beta_0^2 - \bar{y} \beta_0 + \text{const}$$

$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$

$\beta_0^* = \text{argmin}_{\beta_0}$   
 $\beta_0^* = \bar{y}$



$$\frac{1}{2N} \sum_{i=1}^N (y_i - x_i^T \beta)^2 = \frac{1}{2} \beta^T H \beta - \underline{b}^T \beta$$

$H \beta = \underline{b}$   
 $\beta = H^{-1} \underline{b}$

$H = \frac{1}{N} \tilde{X}^T \tilde{X}$ ,  $\underline{b} = \tilde{X}^T \underline{y}$

$\frac{1}{N} \sum_{i=1}^N \tilde{x}_i \tilde{x}_i^T$ ,  $\sum_{i=1}^N \tilde{x}_i y_i$

location

"well-conditioned"  
 Isotropic + Diagonal

$H = I \Rightarrow \beta^T H \beta = \sum_{j=1}^D \beta_j^2$

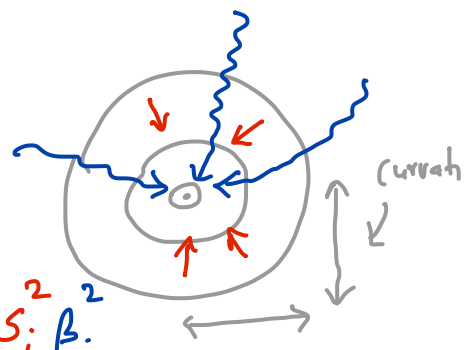
$\beta_1^2 + \beta_2^2$

$\sum_{i=1}^{10,000} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \\ \cdot \end{bmatrix}^T = \tilde{X}^T \tilde{X} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$

ill-conditioned  
 Diagonal

$H = \begin{bmatrix} s_1 & 0 \\ 0 & \dots & 0 \\ 0 & \dots & s_D \end{bmatrix} \Rightarrow \beta^T H \beta = \sum_{j=1}^D s_j^2 \beta_j^2$

off-diagonal



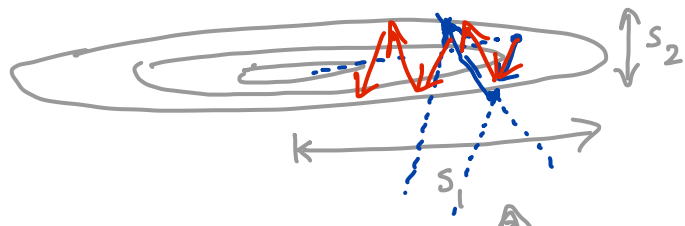
$\begin{bmatrix} s_1 & 0 \\ 0 & s_2 \end{bmatrix} \beta^{(k)} = \begin{bmatrix} s_1 \beta_1^{(k)} \\ s_2 \beta_2^{(k)} \end{bmatrix}$

$\tilde{X}^T \tilde{X} = \begin{bmatrix} s_1^2 & 0 \\ 0 & s_2^2 \end{bmatrix}$

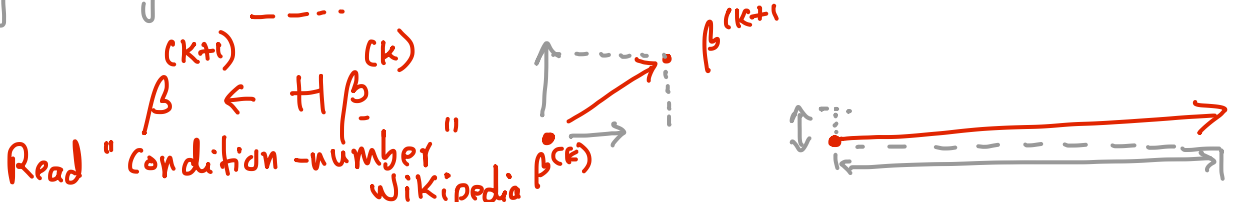
$s_1^2 \beta_1^2 + s_2^2 \beta_2^2$

$s_1^2 \gg s_2^2$

Ratio:  $\frac{s_1}{s_2}$  is really large.



$s_j$ : largest "move" one can make in  $j$ 'th "direction"



General:

$$\tilde{X}^T \tilde{X} \xrightarrow{\text{diagonalization}} U$$

$D \times D$

$$\begin{bmatrix} \uparrow & \uparrow \\ u_1 & \dots & u_k \\ \downarrow & & \downarrow \end{bmatrix} \begin{bmatrix} s_1^2 & & & \\ & s_2^2 & & \\ & & \dots & \\ 0 & & & s_k^2 \end{bmatrix} \begin{bmatrix} \leftarrow & u_i^T & \rightarrow \\ \vdots \\ \leftarrow & u_k^T & \rightarrow \end{bmatrix}$$

$D \times k$        $k \times D$

singular values

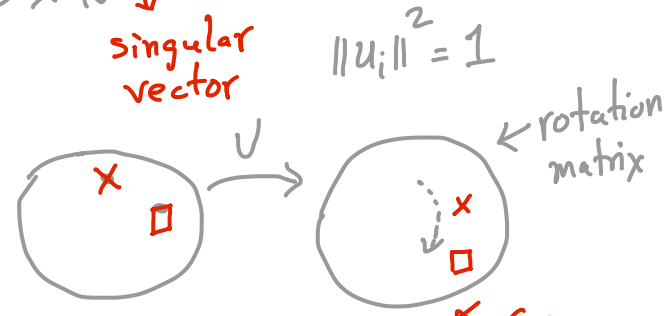
$$\tilde{X}^T = U S V^T \quad (\text{SVD})$$

$D \times N \quad D \times K \quad K \times K \quad K \times N$

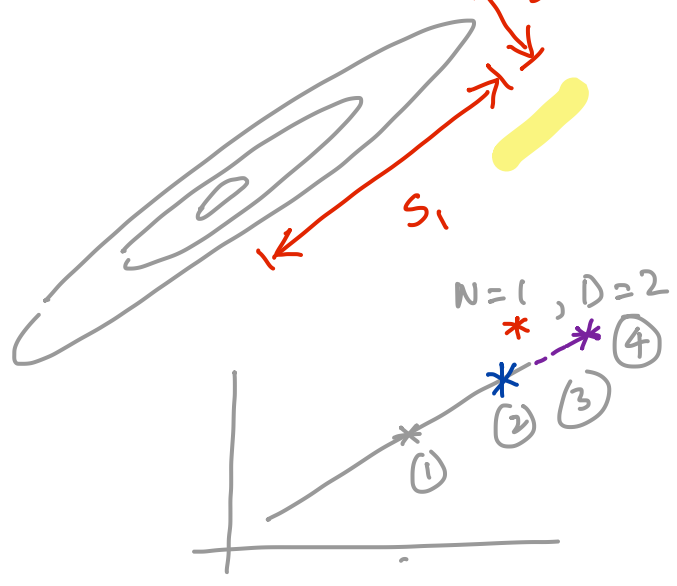
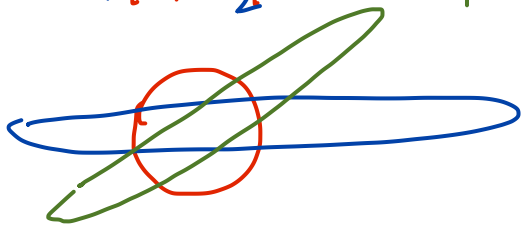
$$V^T V = I_K$$

$$\tilde{X}^T \tilde{X} = (U S V^T) (V S U^T) = U S^2 U^T$$

Read Wikipedia



$$s_1^2 \beta_1^2 + s_2^2 \beta_2^2 + \beta_1 \beta_2$$



Condition number

$$\frac{s_{\max}}{s_{\min}}$$

minimum singular value is 0

Learning rate:

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \alpha^{(k)} X^T (X \beta^{(k)} - y)$$

$\alpha^{(k)} (X^T X) \beta^{(k)}$

effective learning rate (2)

$$\tilde{X}^T = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

$2 \times 1 \quad D \quad N \quad 2 \times 1 \quad 1 \times 1 \quad 1 \times 1$

$s_{\max} > 0$        $K=1$   
 $s_{\min} = 0$        $D=2$

different dimensions "move" differently!  
 $\alpha$  is constrained by  $s_{\max}$  &  $s_{\min}$  (3)

$$\tilde{X}^T = \begin{bmatrix} \cdot & \cdot \\ \cdot & \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \square & \square \end{bmatrix} \begin{bmatrix} \cdot \\ \cdot \end{bmatrix}$$

$2 \times 2$        $2 \times 1 \quad 1 \times 1 \quad 1 \times 1$

$\tilde{X}^T = \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} = \begin{bmatrix} \cdot \\ \cdot \end{bmatrix} \begin{bmatrix} \square & \square \end{bmatrix} \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix}$

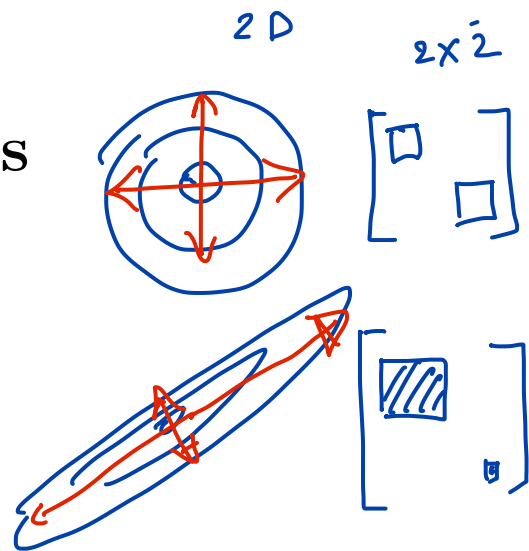
$2 \times 3$        $2 \times 1 \quad 1 \times 1 \quad 1 \times 1$

3rd

## Invertibility and uniqueness

The Gram matrix  $\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}$  is invertible iff  $\tilde{\mathbf{X}}$  has full column rank.

Proof: Assume  $N > D$ . The fundamental theorem of linear algebra states that the dimensionality of null space is zero for full column rank. This implies that the Gram matrix is positive definite, which implies invertibility.



## Rank deficiency and ill-conditioning

Unfortunately,  $\tilde{\mathbf{X}}$  could often be rank deficient in practice, e.g. when  $D > N$ , or when the columns  $\bar{\mathbf{x}}_d$  are (nearly) collinear. In the later case, the matrix is ill-conditioned, leading to numerical issues.

## Summary of linear regression

We have studied three methods:

1. Grid search
2. (Stochastic) gradient descent
3. Least squares



$$\left(\frac{x_i - \bar{x}_i}{\|x_i - \bar{x}_i\|}\right) \quad \frac{x_i}{\|x_i\|}$$

# Additional Notes

Normalization:  $(X^T X)_{ij} = \sum_{i=1}^N x_{i1} x_{ij}$

$D \times 1$   $N \times 1$

$X^T X = \begin{bmatrix} \|x_1\|^2 & (x_1^T x_2) & \dots \\ \dots & \dots & \dots \\ \dots & \dots & \|x_D\|^2 \end{bmatrix}$

$\sum_{i=1}^N x_i x_i^T$

$[1, x_{n1}, x_{n2}, \dots, x_{nd}]^T \quad D \times 1$

## Implementation

There are many ways to implement matrix inversion, but using QR decomposition is one of the most robust ways. Matlab's backslash operator implements this (and much more) in just one line.

```

1 beta = inv(X'*X) * (X'*y)
2 beta = pinv(X'*X) * (X'*y)
3 beta = (X'*X) \ (X'*y)

```

For robust implementation, see Sec. 7.5.2 of Kevin Murphy's book.

## To do

$K = D \text{ or } N$  (none of  $S_j$  are 0)

1. Revise linear algebra to understand why  $\tilde{X}$  needs to have full rank. Read the Wikipedia page on rank of a matrix.  $S_j \neq 0 \forall j$
2. For details on the geometrical interpretation, see Bishop 3.1.2. However, better to read this after the lecture on "basis-function expansion". Also, note that notation in the book is different. This might make the reading difficult.  $\downarrow$  infinite sol<sup>n</sup>
3. Understand matrix inversion robust implementation and play with it using the code for labs. Read Kevin Murphy's section 7.5.2 for details.
4. Understand ill-conditioning. Reading about the "condition number" in Wikipedia will help. Also, understanding SVD is essential. Here is another link provided by Dana Kianfar (EPFL) <http://www.cs.uleth.ca/~holzmann/notes/illconditioned.pdf>.
5. Work out the computational complexity of least-squares (use the Wikipedia page on computational complexity).

Read this note

# 7 Challenge: Overfitting

## Motivation

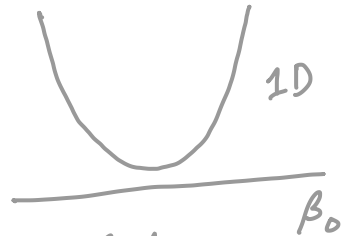
Linear model can be easily modified to obtain more powerful non-linear model. We can use basis function expansion to get a non-linear regression model, and then use a sequence of these models to construct a deep model.

Consider simple linear regression. Given one-dimensional input  $x_n$ , we can generate a polynomial basis.

$$\phi(x_n) = [1, x_n, x_n^2, x_n^3, \dots, x_n^M]$$

Then we fit a linear model using the original *and* the generated features:

$$y_n \approx \beta_0 + \beta_1 x_n + \beta_2 x_n^2 + \dots + \beta_M x_n^M$$



$$y_n \approx \beta_0 \cdot 1$$

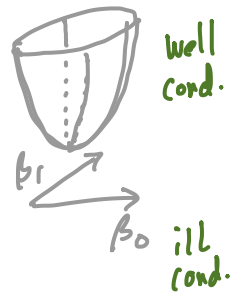
$$y_n \approx \beta_0 \cdot 1 + \beta_1 \cdot x_n$$

Suppose  $x_n = c$ , a constant. We expect in 2D

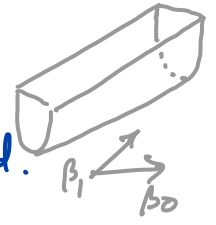
$$y_n \approx \beta_0 \cdot 1 + \beta_1 c$$

$$= \beta_0^{new} \cdot 1$$

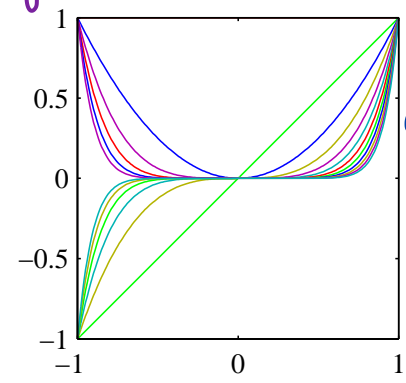
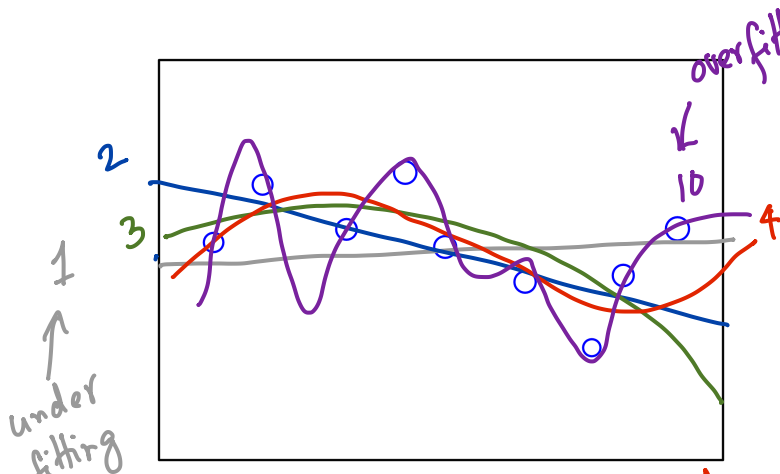
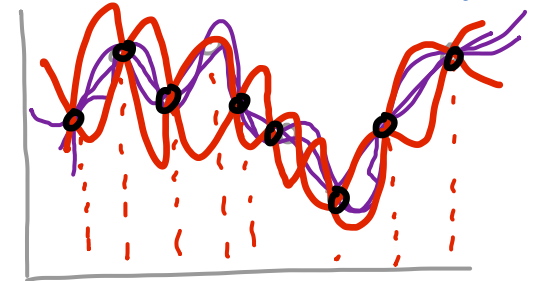
where  $\beta_0^{new} = \beta_0 + c\beta_1$



Adding more data "may" make the loss ill-conditioned.

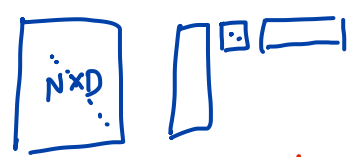


interplay bet<sup>n</sup>  $N$  &  $D$ , but what matters is  $k$  (rank) of  $X$ .



Infinite sol<sup>n</sup>  
 $\downarrow$   
 e.g. some parameters (or their combinations)  
 In linear model whenever  $k < \min(N, D)$   
 $\downarrow$   
 there are some zero singular values?  
 $S_j = 0, j > k+1 \dots$

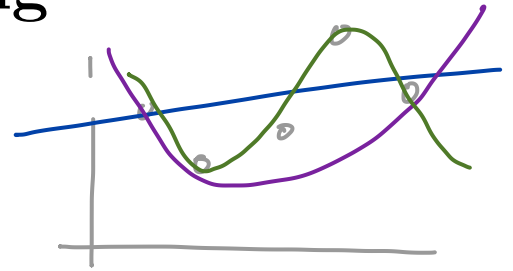
$N$  is fixed,  $D$  is increased  
 Is that good?



Simplest fix: Process data, feature extraction/engineering.

# Overfitting and Underfitting

**Overfitting** is fitting the noise in addition to the signal. **Underfitting** is not fitting the signal well. In reality, it is very difficult to be able to tell the signal from the noise.

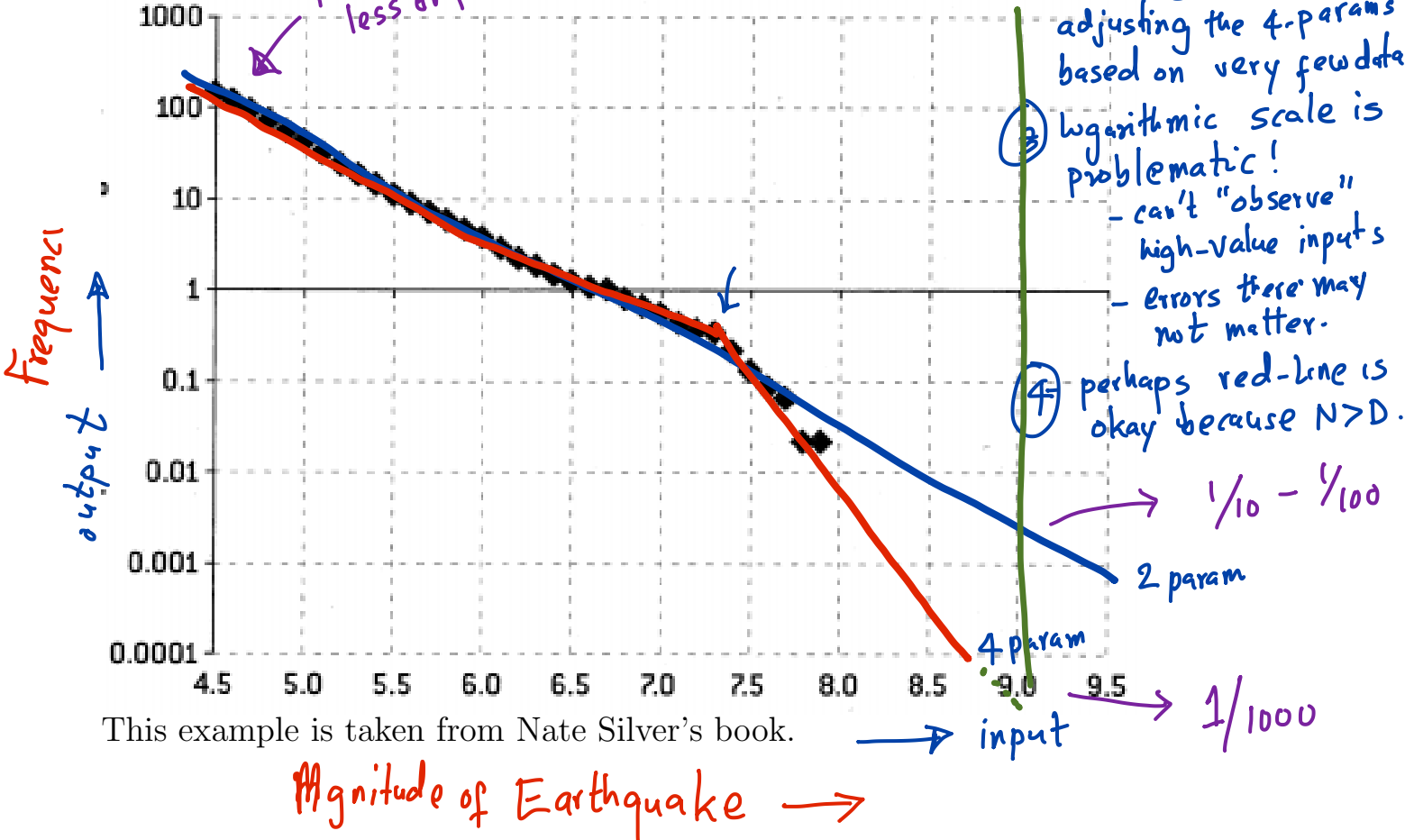


## Which is a better fit?

Try a real situation. Below, y-axis is the frequency of an event and x-axis is the magnitude. It is clear that as magnitude increases, frequency decreases.

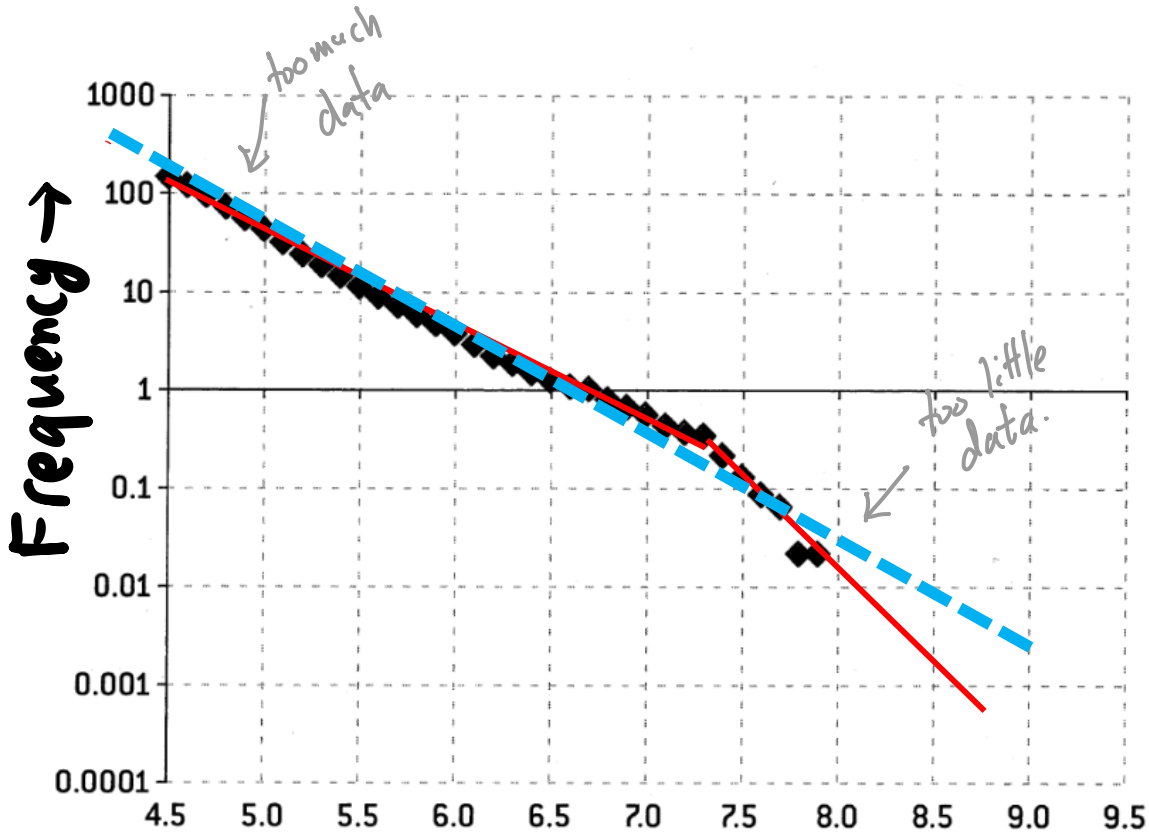
Tohoku

prediction is less important

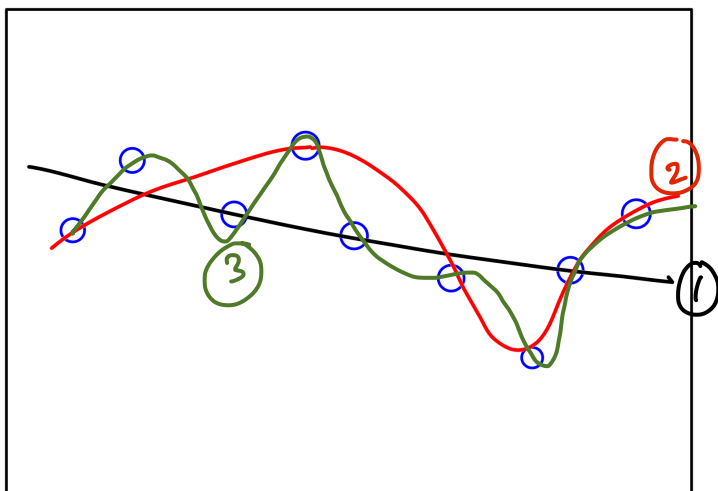


This example is taken from Nate Silver's book.

Which model is a better fit? blue or red?



Another example: Which model is a better fit? black or red? Data is denoted by circle.



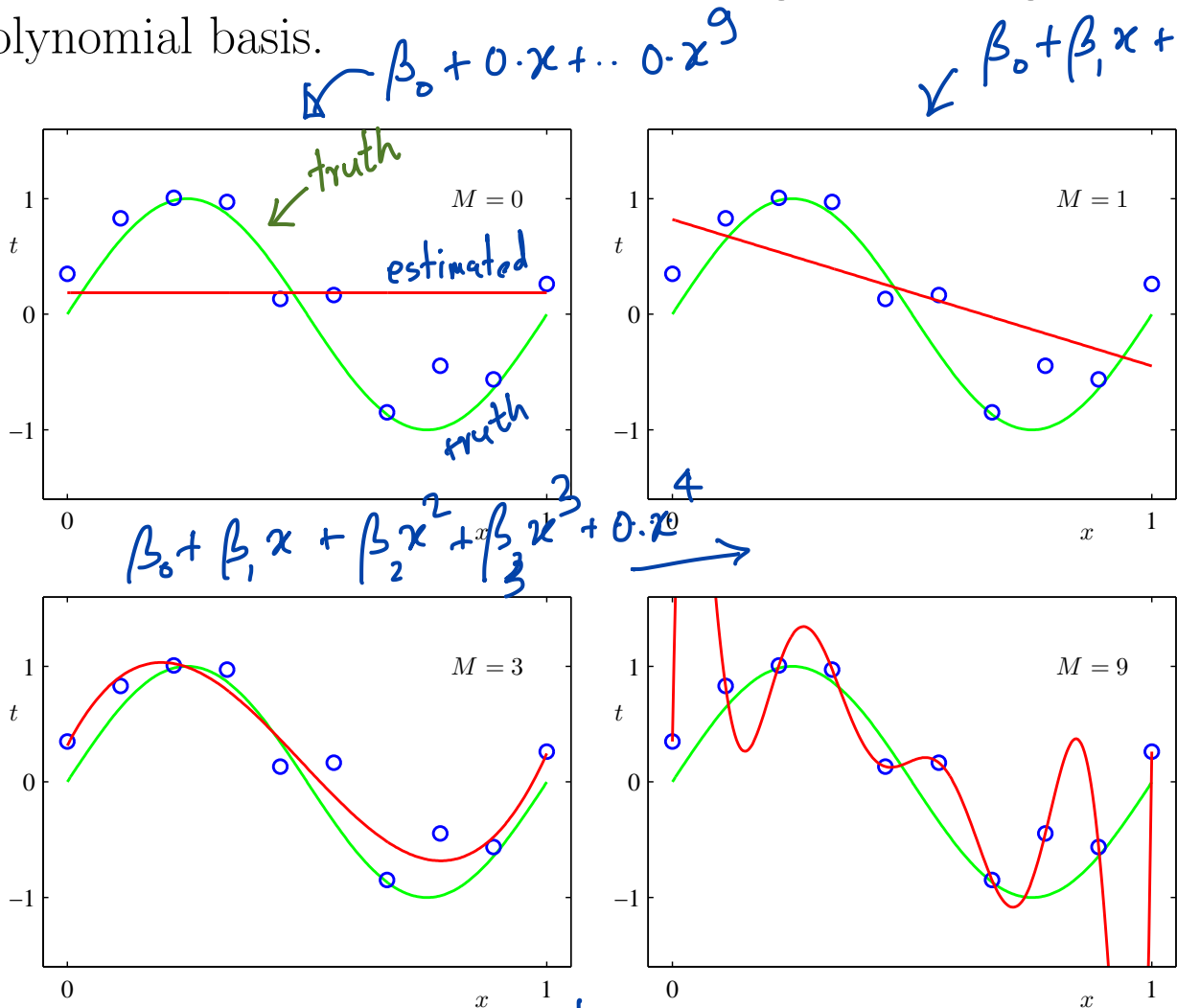
MSE over "data in hand"

- (1) 0.72
- (2) 0.2
- (3) 0

This is not the right criteria to look at.

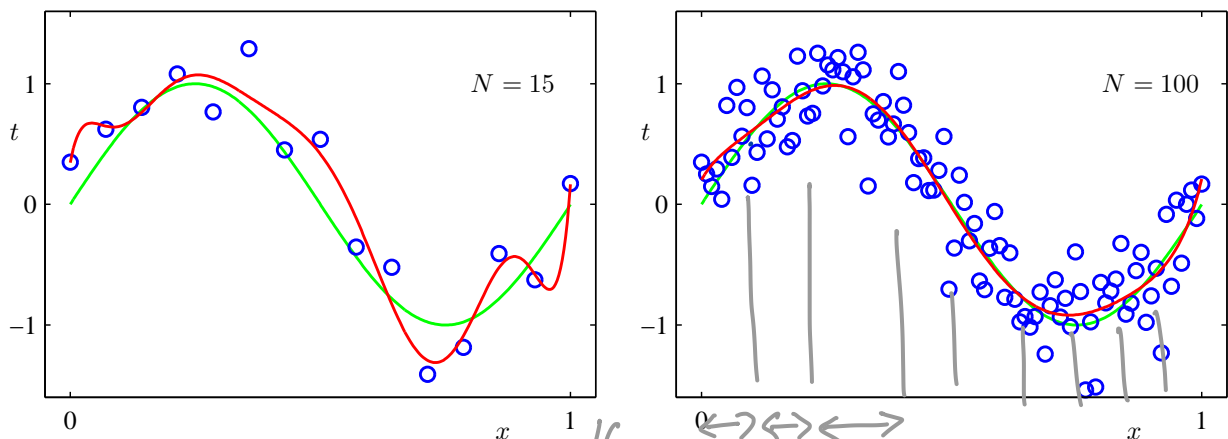
# Complex models overfit easily (see Bishop)

Circles are data points, green line is the truth & red line is the model fit.  $M$  is the maximum degree in the generated polynomial basis.



When we can collect the data,

If you increase the amount of data, overfitting *might* reduce.



$4 \times 4 = 16 \leftarrow \boxed{256}^{16}$

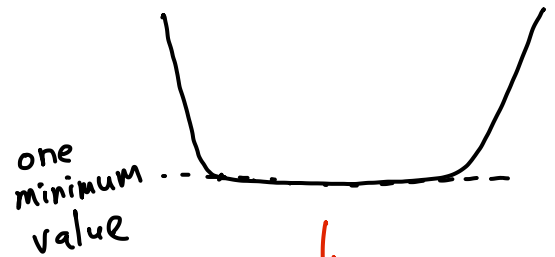
# Occam's razor

One solution is dictated by Occam's razor which states that "Simpler models are better – in absence of certainty."

Sometimes, if you increase the amount of data, you might reduce overfitting. But, when unsure, choose a simple model over a complicated one.

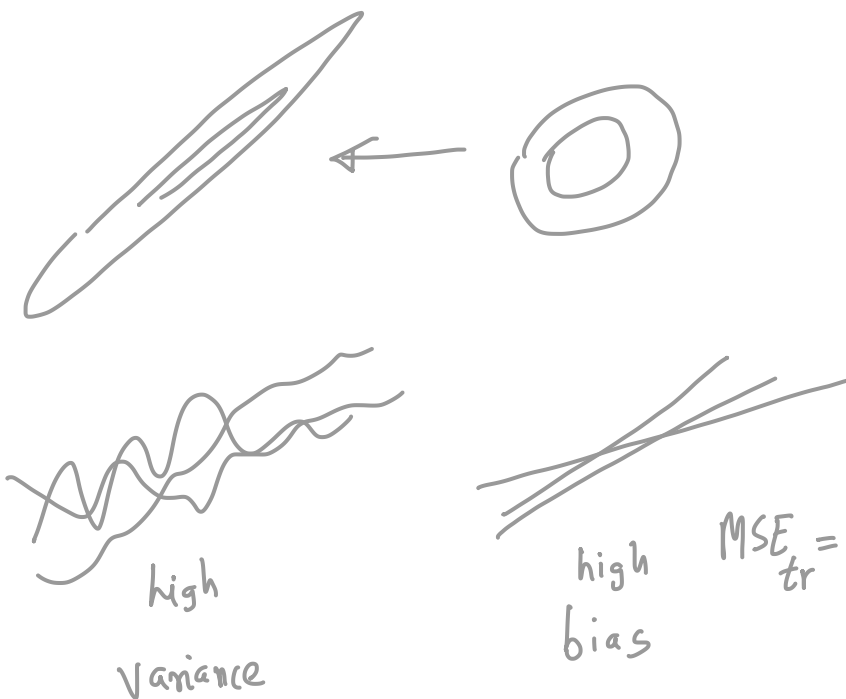
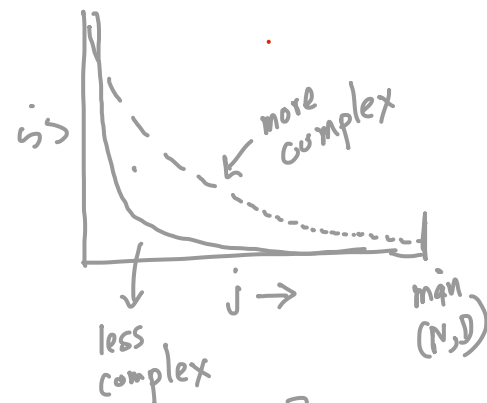
Complexity

Linear model  $L$   
+  
MSE  
+  
any algorithm



Rank of  $X^T X$   
 $= U S^2 U^T$

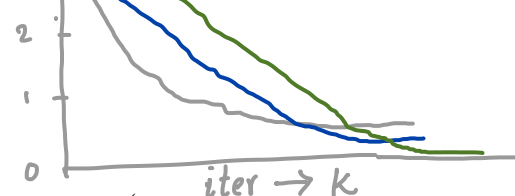
(in decreasing order)



$$MSE_{tr} = \sum_{i \in \mathcal{D}_{tr}} \left( y_i - \tilde{x}_i^T \beta \right)^2 \quad \forall i \text{ in data}$$

all data in hand

model at k'th iteration



## Additional Notes

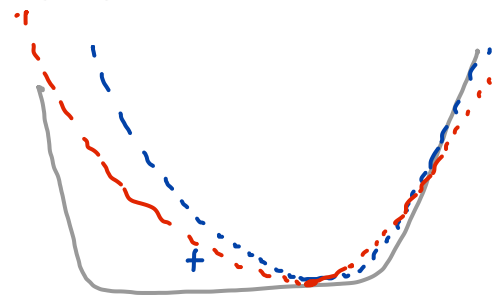
Read about overfitting in the paper by Pedro Domingos (section 3 and 5 of "A few useful things to know about machine learning"). You can also read Nate Silver's book on "The signal and the noise" (the earthquake example is taken from this book).

Complex model  $\rightarrow$  ill-conditioning  $\rightarrow$  Overfitting

## 8 Solutions: Regularization

### What is regularization?

Through regularization, we can penalize complex models and favor simpler ones:



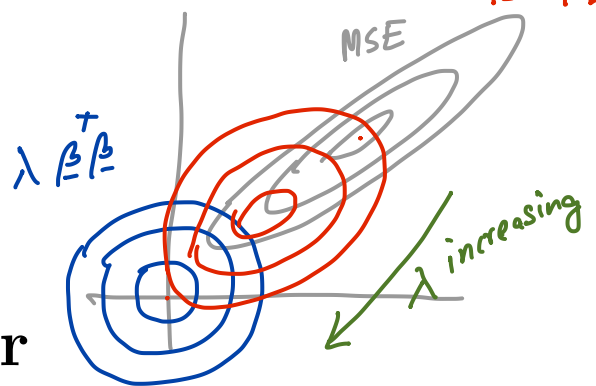
$$\min_{\beta} \text{MSE } \mathcal{L}(\beta) + \frac{\lambda}{2N} \sum_{j=1}^M \beta_j^2 \quad \text{Regularizer} \quad \sum_{i=1}^N (y_i - \beta_0 - \beta_1 x_{ii})^2 + \lambda (\beta_0^2 + \beta_1^2) \quad \text{(A)}$$

### Ridge Regression

The second term is a regularizer (with  $\lambda > 0$ ). The main point here is that an input variable weighted by a small  $\beta_j$  will have less influence on the output.

Q1: What happens,  $\lambda \rightarrow \infty$   $\beta^* = 0$

Q2: " " ,  $\lambda \rightarrow 0$   $\beta^* = \beta_{LS}$



### Regularization Parameter

The parameter  $\lambda$  can be tuned to reduce overfitting. But, how do you choose  $\lambda$ ? Write the SVD of  $X^T = USV^T$

Q: Why "simple" features are favoured?

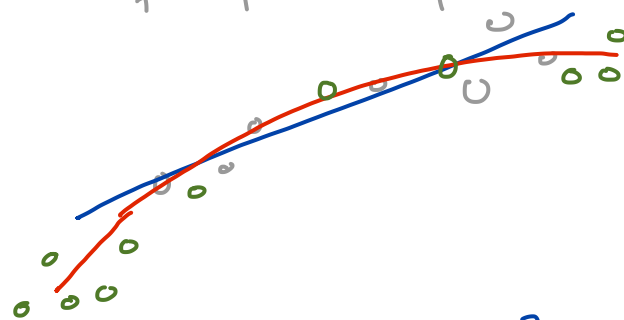
Q: Write an expression for the minimizer of (A)

$$\beta^*(\lambda) = (X^T X + \lambda I)^{-1} X^T Y$$

Q: How does the conditioning of the problem improved?

### The generalization error

The generalization error of a learning method is the expected prediction error for "unseen" data, i.e. mistakes made on the data that we are going to see in the future. This quantifies how well the method generalizes.



What is "unseen" data?

Tricky to define ...  
nevertheless ....

$$\begin{aligned}
 & X^T X + \lambda I \\
 &= U S V^T V S U^T + \lambda I_D \\
 &= U \underbrace{S^2}_{I_K} U^T + \lambda U U^T \\
 &= U (S^2 + \lambda I) U^T
 \end{aligned}$$

Q: Write an expression for the minimizer of  $\|A^{-1} \tau\|$

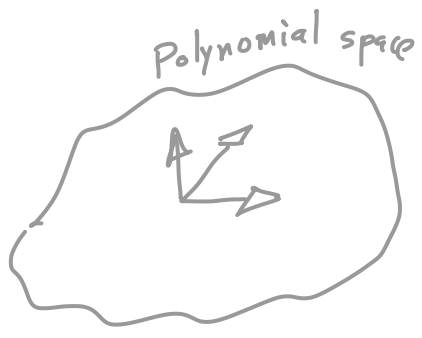
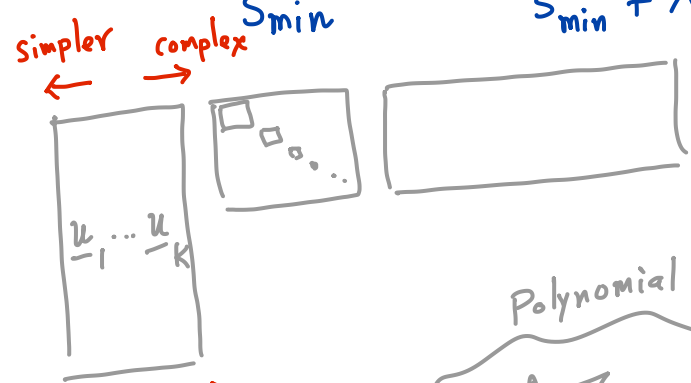
$$\beta^*(\lambda) = (X^T X + \lambda I)^{-1} X^T y$$

Q: How does the conditioning of the problem improved?

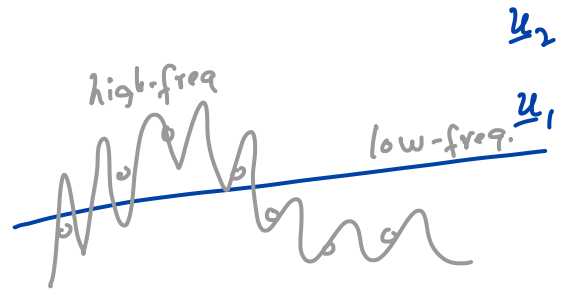
Write the SVD of  $X^T = U S V^T$

Q: Why "simpler" features are favoured

Cond. num.  $\frac{S_{max}^2}{S_{min}^2} \rightarrow \frac{S_{max}^2 + \lambda}{S_{min}^2 + \lambda}$



as  $\lambda \rightarrow \infty$ , cond. num.  $\rightarrow 1$



$$X^T = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_{1,2} & x_{2,2} & \dots & x_{N,2} \\ \vdots & \vdots & \dots & \vdots \\ x_{1,D} & x_{2,D} & \dots & x_{N,D} \end{bmatrix}$$

$D \times N$



# Simulating the future

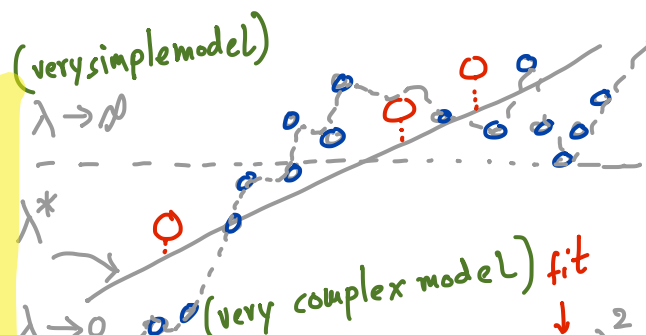
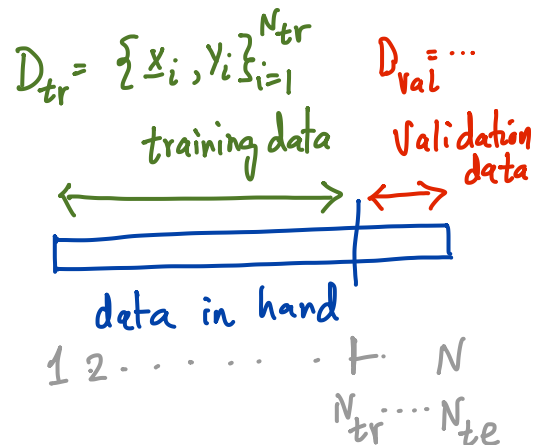
Ideally, we should choose  $\lambda$  to minimize the mistakes that will be made in the future. Obviously, we do not have the future data, but we can always *simulate the future* using the data in hand.

# Splitting the data

For this purpose, we split the data into train and validation sets, e.g. 80% as training data and 20% as validation data. We pretend that the validation set is the future data. We fit our model on the training set and compute a prediction-error on the validation set. This gives us an *estimate* of the generalization error (one instant of the future).

$$\beta_0 + \beta_1 x + \dots + \beta_{100} x^{100}$$

$\lambda \rightarrow \infty, \beta_j \rightarrow 0$



$$MSE_{tr}(\beta) = \sum_{i \in D_{tr}} (y_i - \tilde{x}_i^T \beta)^2$$

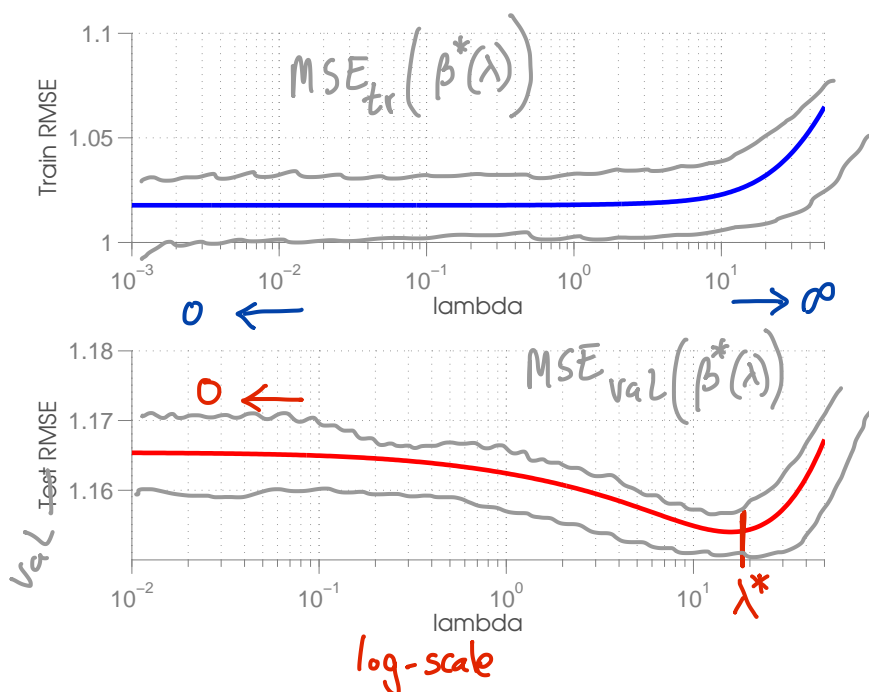
$$MSE_{val}(\beta) = \sum_{i \in D_{val}} (y_i - \tilde{x}_i^T \beta)^2$$

$$\beta^*(\lambda) = \min_{\beta} MSE_{tr}(\beta) + \frac{\lambda}{2} \beta^T \beta$$

Q: How should we choose the split?

A: randomly? but we should be careful.

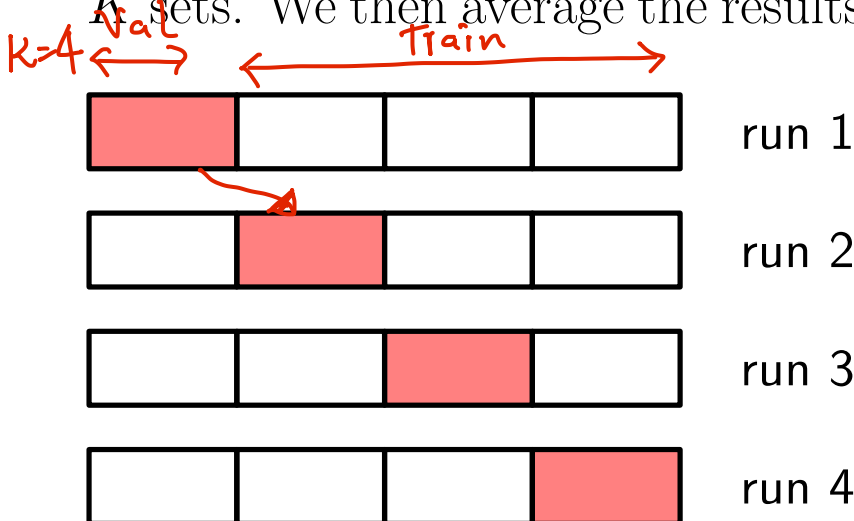
It should be "representative" of the problem/data.



# Cross-validation

Random splitting (aka bootstrap) is not an efficient method.

K-fold cross-validation allows us to do this efficiently. We randomly partition the data into  $K$  groups. We train on  $K - 1$  groups and test on the remaining group. We repeat this until we have tested on all  $K$  sets. We then average the results.



Cross-validation returns an estimate of the *generalization error*.

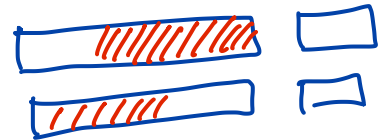
## Additional Notes

Details on cross-validation are in Chapter 7 in the book by Hastie, Tibshirani, and Friedman (HTF). You can also read about bootstrap in Section 7.11 in HTF book. This method is related to random splitting and is a very popular method.

for  $s = 1, \dots, 100$   
 shuffle the data  
 pick the first  $x\%$  as train  
 fit  
 evaluate  $MSE_{val}(\beta_s^*(\lambda))$

end

Random shuffling is not efficient!



Validation sets are disjoint!  
 Estimates of  $\text{valError}$  is a bit better (low variance) but it's biased.

shuffle and  
 \* Split in  $K$  subsets

\*  $\lambda = [10^{-5} \ 10^{-4} \ \dots \ 1 \ 100 \ 10^3]$

\* for all  $\lambda$

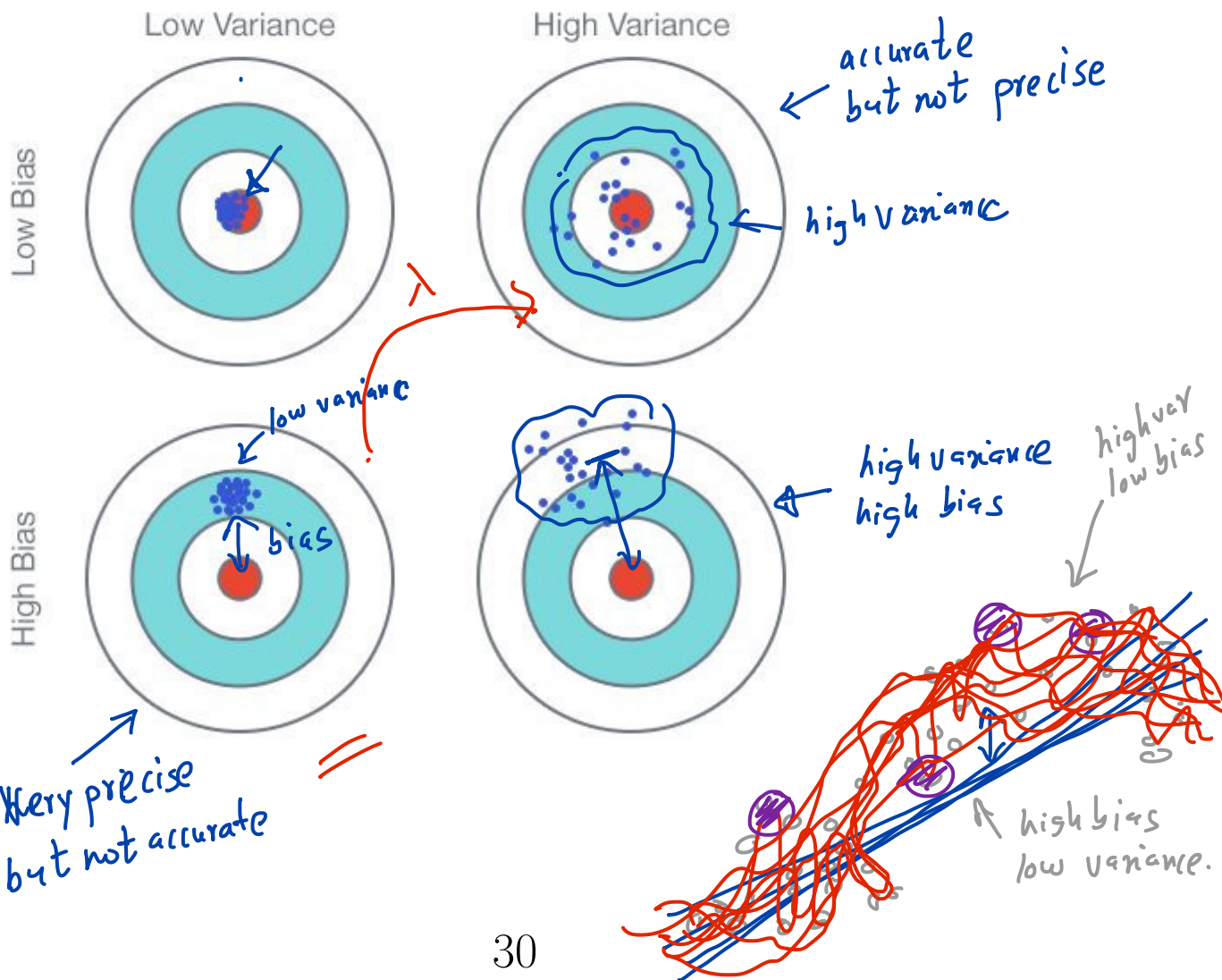
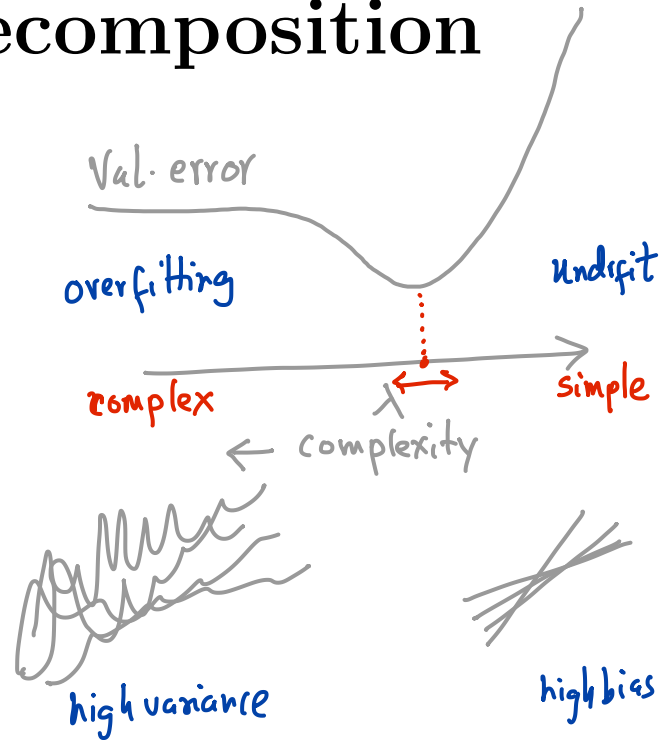
for  $k = 1 : K$   
 Fit model Train  
 test on the validation  
 $MSE_{val}^{\lambda}[k]$   
 $MSE_{val}^{\lambda} = \text{mean over } k=1:K [MSE_{val}^{\lambda}[k]]$   
 end

# 9 Bias-Variance Decomposition

## What is bias-variance?

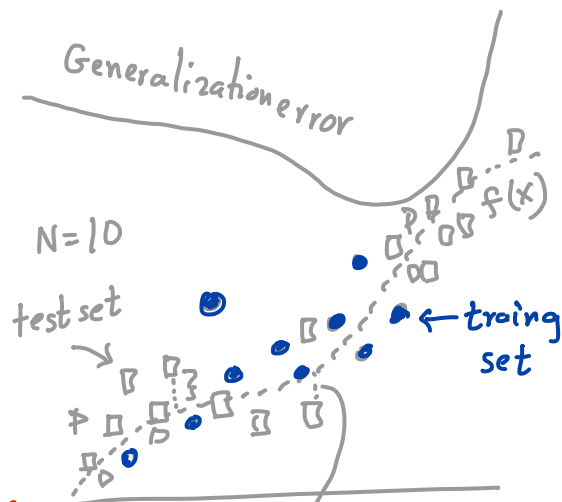
One natural question is how does the test error vary wrt  $\lambda$ ? When  $\lambda$  is high, the model underfits, while when  $\lambda$  is small, the model overfits. Therefore, a good value is somewhere in between.

Bias-variance decomposition explains the shape of this curve.



# Generalization error

Given training data  $D_{tr}$  of size  $N$ , we would like to estimate the expected error made in future prediction. This error is the generalization error. Below is a definition suppose that we have infinite test data  $D_{te}$ ,



$$teErr(D_{tr}) := \mathbb{E}_{D_{te}} [\{y - f(\mathbf{x})\}^2] = \mathbb{E}_{y, \mathbf{x} \sim D_{te}} [(y - \beta_*^T \tilde{\mathbf{x}})^2]$$

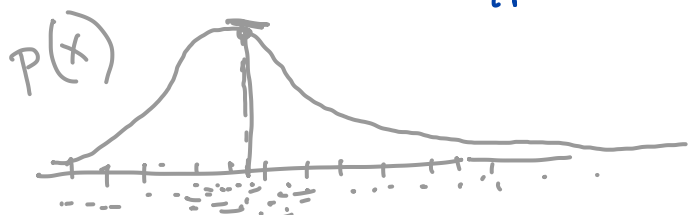
↑ training data
↑ training data
↑ test data
↑ best linear model obtained by

Generalization error is different from the training error which measures how well you fit the data.

$$trErr(D_{tr}) := \sum_{n=1}^N [\{y_n - f(\mathbf{x}_n)\}^2]$$

$$\beta_*^{(D_{tr})} = \arg \min_{\beta} \frac{1}{N} \sum_{\mathbf{x}_n, y_n \sim D_{tr}} [(y_n - \beta^T \tilde{\mathbf{x}}_n)^2]$$

↑ function of  $D_{tr}$   
 ↑ function of training data  $D_{tr}$



$$\lim_{n \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N x_n \rightarrow \mathbb{E}(x) = \int x P(x) dx$$

mean

$$\text{Mean-Squared error} = \frac{1}{N} \sum_{n=1}^N (y_n - f(x_n))^2 \quad D_{te} = \{y_n, x_n\}_{n=1}^{\infty}$$

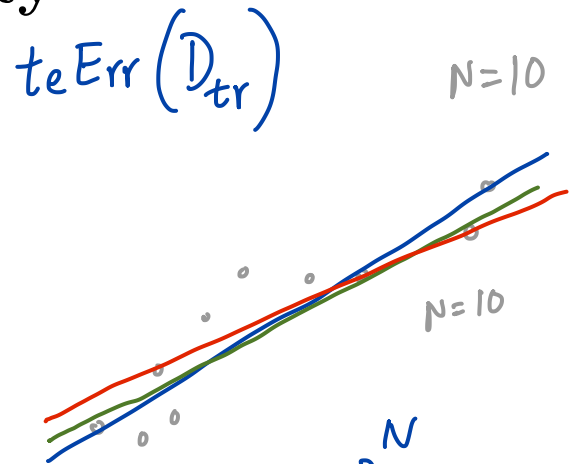
$$\mathbb{E}_{D_{te}} [(y - f(x))^2] = \int (y - f(x))^2 P(y, x) dy dx$$

↓  $N \rightarrow \infty$

# Errors vs model complexity

As we increase the model complexity, how do these errors vary? The blue line shows training error for a dataset with  $N = 50$ , while the red line shows the generalization error for that dataset.

Simple model have high train and generalization error since they have a high **bias**, while complex model have low train but high generalization error because they have high **variance**.



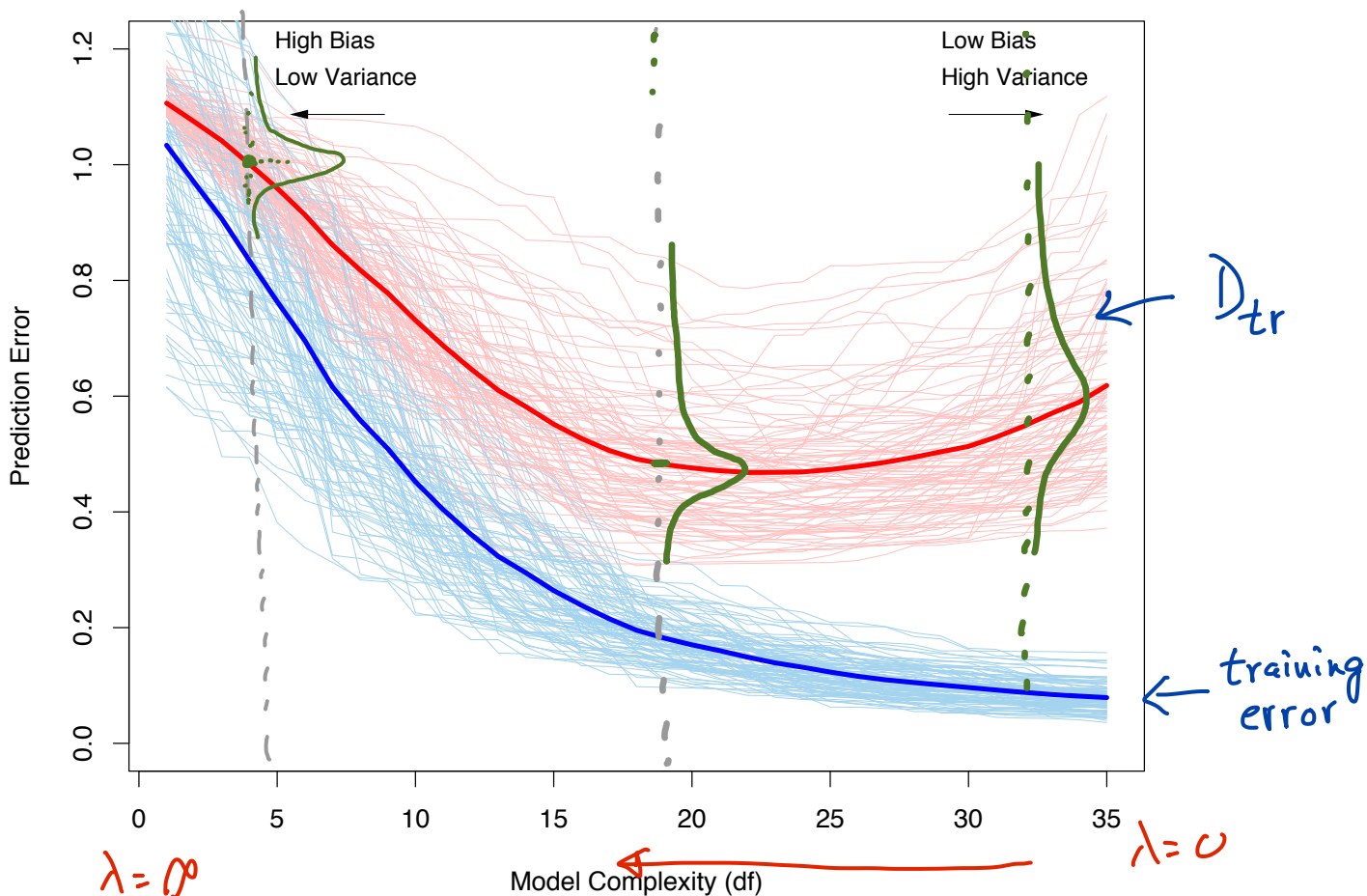
$$D_{tr} = \{x_n, y_n\}_{n=1}^N$$

where  $x_n, y_n \sim P(x, y)$

$$\beta_*(D_{tr}; \lambda)$$

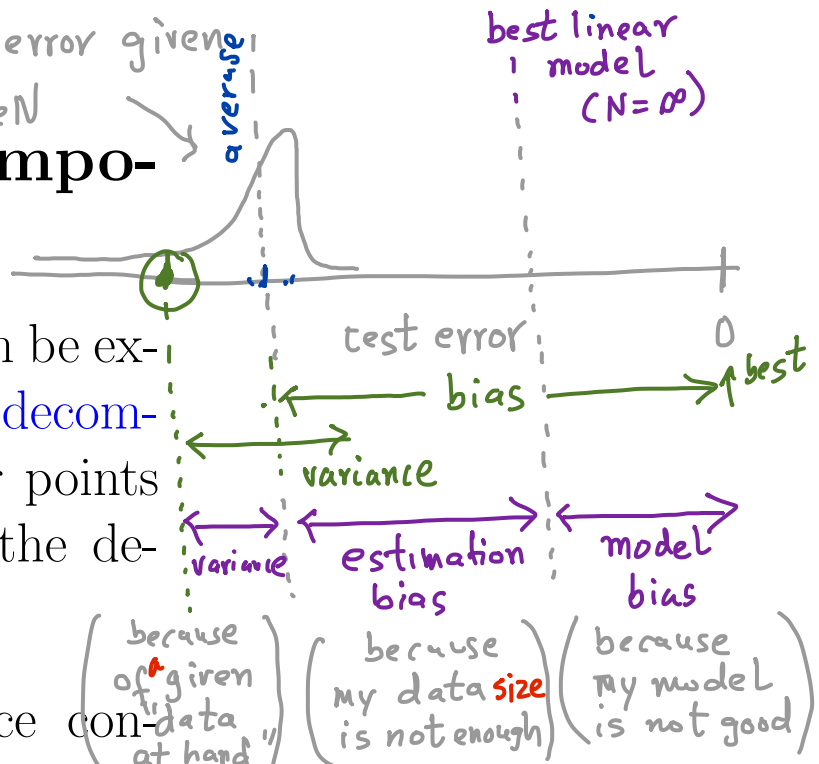
$\downarrow$

$$teErr(D_{tr}; \lambda)$$



# Bias-variance decomposition

Distribution of test error given  $D_L \sim P(x, y)$  of size  $N$



The shape of these curves can be explained using **bias-variance decomposition**. The following four points can be explained by using the decomposition:

1. both bias and variance contribute to generalization error.

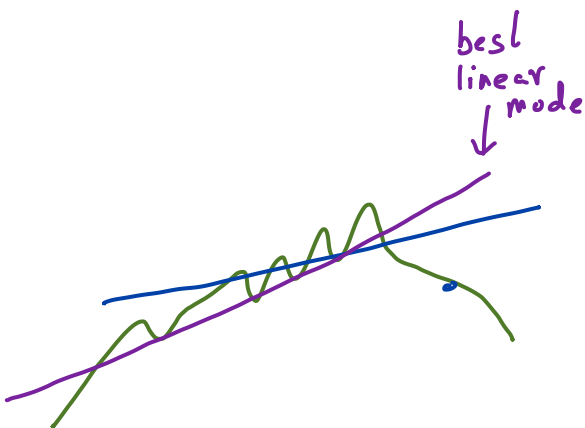
2. For bias, both **model-bias** and **estimation-bias** are important.

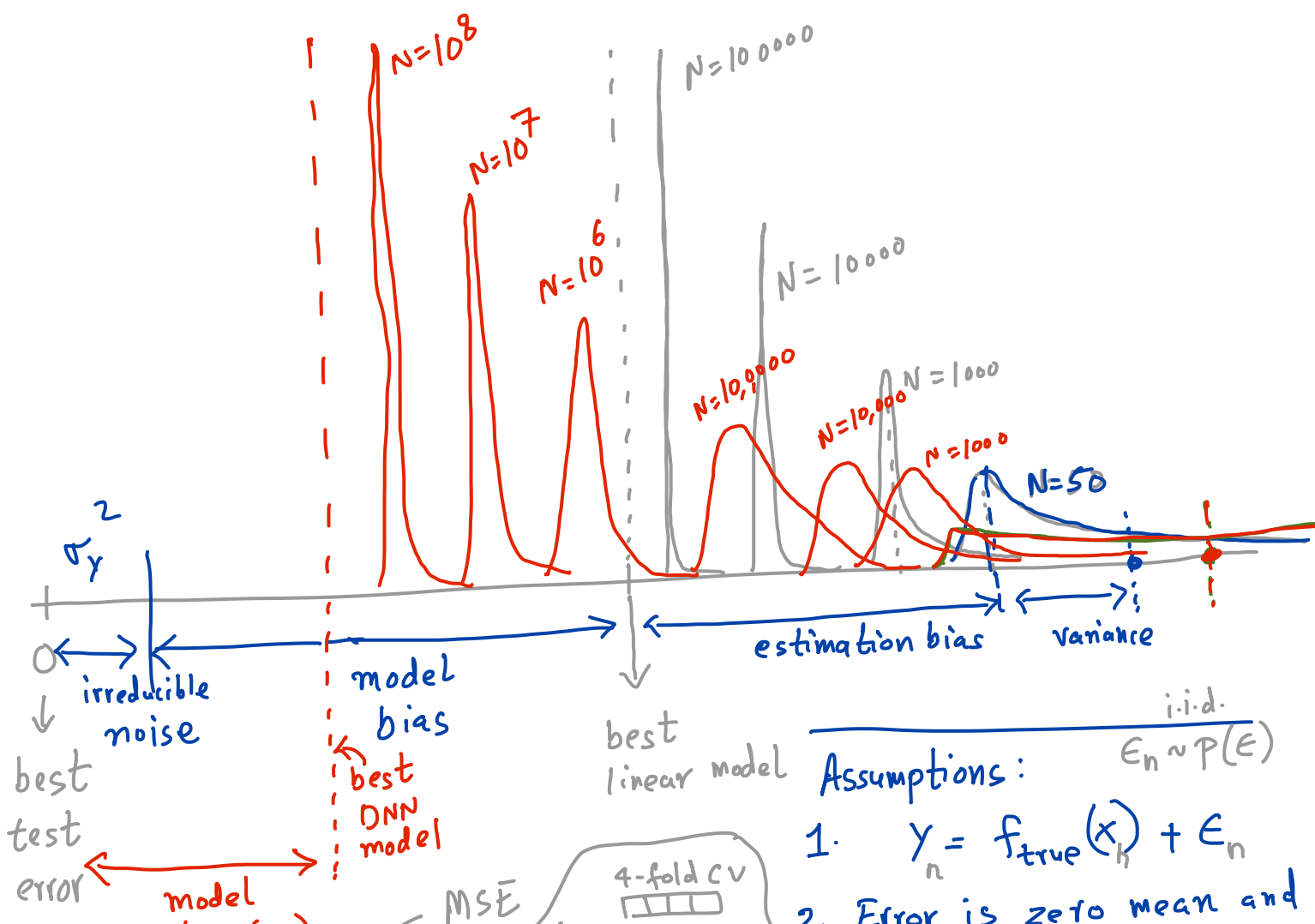
When we increase model complexity, we increase generalization error due to increased variance.

3. Regularization increases estimation bias while reducing variance.

We will use "add-subtract" technique to analyze this:

$$\begin{aligned}
 & y - \beta_*^T \tilde{x} \\
 = & \underbrace{y - f_{\text{true}}(x)} + \underbrace{f_{\text{true}}(x) - \beta_*^T \tilde{x}} \\
 & \underbrace{f_{\text{true}}(x) - \beta_{\text{best}}^T \tilde{x}} + \underbrace{\beta_{\text{best}}^T \tilde{x} - \beta_*^T \tilde{x}} \\
 & \text{model bias} \qquad \text{estimation bias + variance}
 \end{aligned}$$





Assumptions:  $E_n \sim P(E)$  i.i.d.

- $y_n = f_{\text{true}}(x) + \epsilon_n$
- Error is zero mean and variance  $\sigma_y^2$   
(Doesn't need to be Gaussian)
- $x$  is "given", not stochastic

$$\mathbb{E}_{D_{\text{te}}}[\dots] = \mathbb{E}_{P(y)}[\dots]$$

$$\mathbb{E}_{D_{\text{tr}}} \mathbb{E}_{D_{\text{te}}} \left[ (y - f_{\text{est}}(x))^2 \right] \approx \sum_{i=1}^4 \sum_{j \in D_{\text{te}}^{(i)}} (y_j - f(x_j))^2$$

4-fold cv

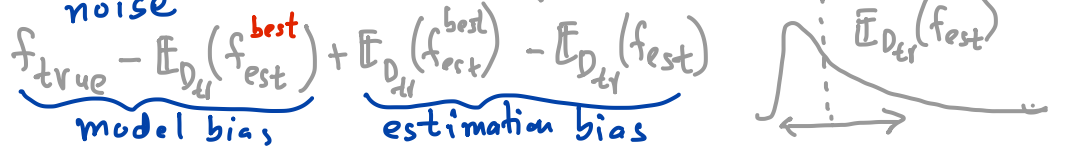
$$= \mathbb{E}_{D_{\text{tr}}} \mathbb{E}_{D_{\text{te}}} \left[ (y - f_{\text{true}} + f_{\text{true}} - f_{\text{est}})^2 \right]$$

$$= \underbrace{\mathbb{E}_{D_{\text{te}}} [(y - f_{\text{true}})^2]}_{\text{irreducible noise}} + \underbrace{\mathbb{E}_{D_{\text{tr}}} [(f_{\text{true}} - f_{\text{est}})^2]}_{\text{bias}} + 2 \underbrace{\mathbb{E}_{D_{\text{te}}} (y - f_{\text{true}})}_{\epsilon} \underbrace{\mathbb{E}_{D_{\text{tr}}} (f_{\text{true}} - f_{\text{est}})}_{\text{estimation bias}}$$

$$= \sigma_y^2 + \mathbb{E}_{D_{\text{tr}}} [(f_{\text{true}} - f_{\text{est}})^2]$$

$$= \sigma_y^2 + \mathbb{E}_{D_{\text{tr}}} \left( f_{\text{true}} - \mathbb{E}_{D_{\text{tr}}}(f_{\text{est}}) + \mathbb{E}_{D_{\text{tr}}}(f_{\text{est}}) - f_{\text{est}} \right)^2$$

$$= \underbrace{\sigma_y^2}_{\text{irreducible noise}} + \underbrace{\left( f_{\text{true}} - \mathbb{E}_{D_{\text{tr}}}(f_{\text{est}}) \right)^2}_{\text{bias}} + \mathbb{E}_{D_{\text{te}}} \left[ \left( f_{\text{est}} - \mathbb{E}_{D_{\text{tr}}}(f_{\text{est}}) \right)^2 \right] + 2 \left( f_{\text{true}} - \mathbb{E}_{D_{\text{tr}}}(f_{\text{est}}) \right) \times \underbrace{\mathbb{E}_{D_{\text{tr}}} \left( f_{\text{est}} - \mathbb{E}_{D_{\text{tr}}}(f_{\text{est}}) \right)}_{=0}$$



# Understanding bias - variance "Trade-off" in Linear model.

(Ridge regression)

$$\beta_{LS}^*(D_{tr}) = \left( \tilde{X}_{tr}^T \tilde{X}_{tr} \right)^{-1} \tilde{X}_{tr}^T Y_{tr}$$

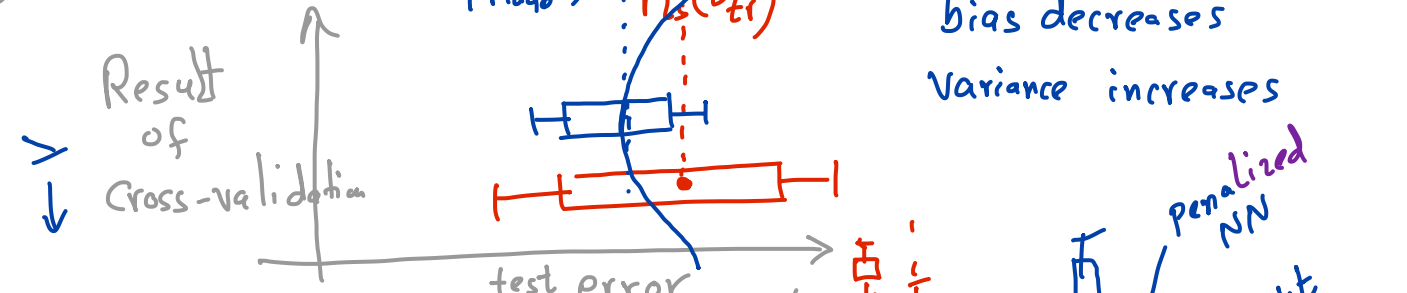
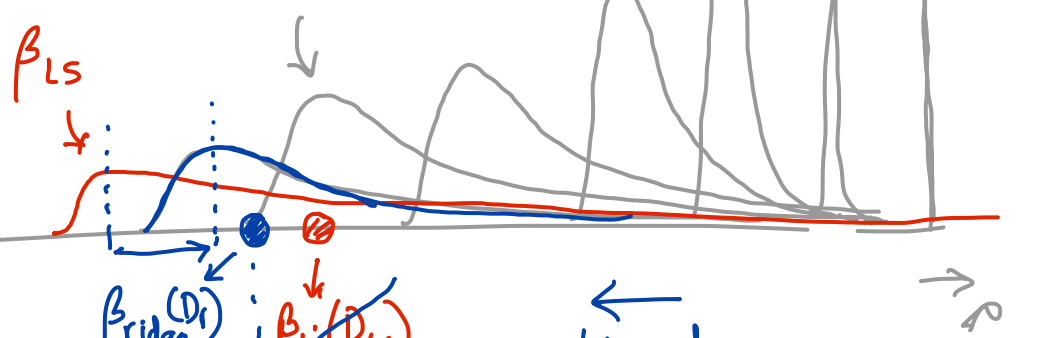
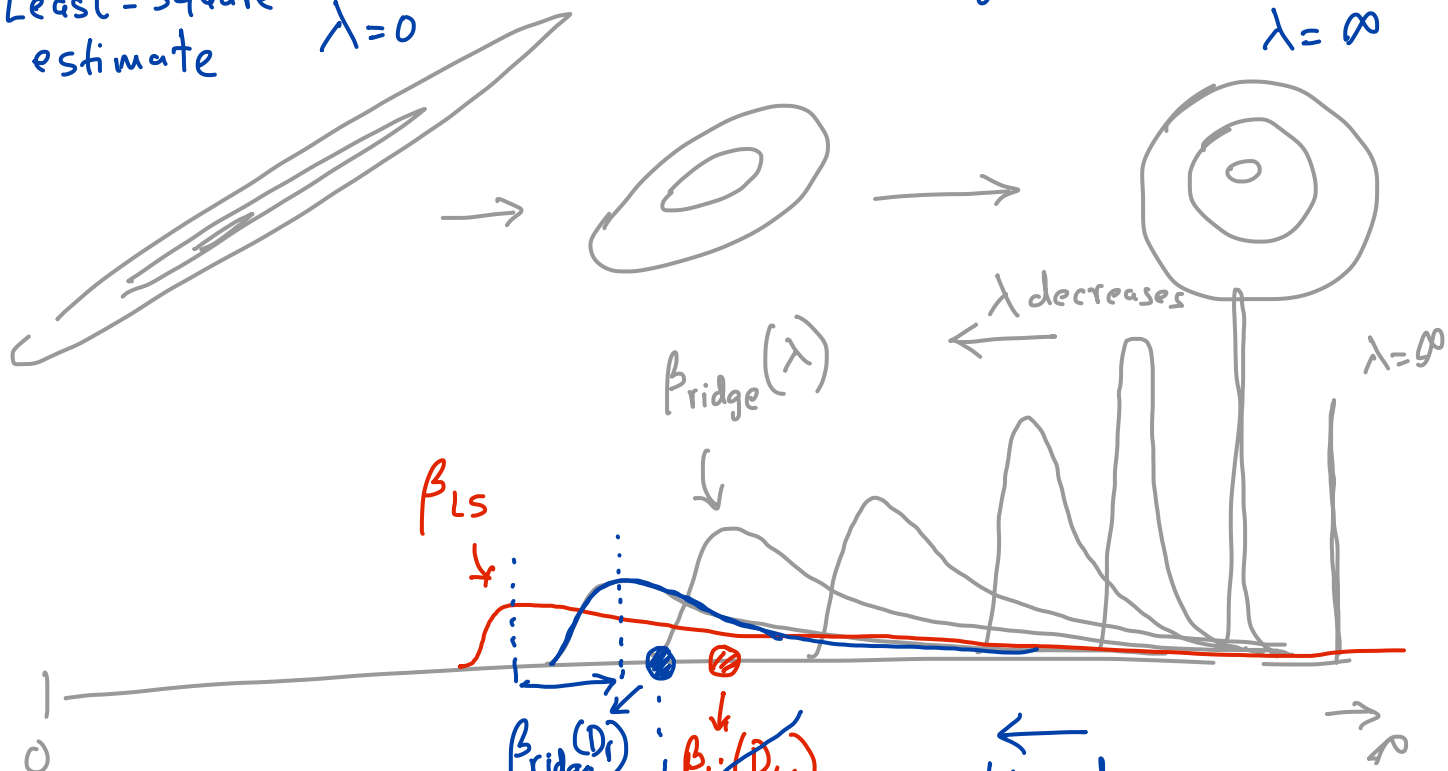
Least-square estimate

$\lambda = 0$

$$\beta_{Ridge}^*(D_{tr}, \lambda) = \left( \tilde{X}_{tr}^T \tilde{X}_{tr} + \lambda I \right)^{-1} \tilde{X}_{tr}^T Y_{tr}$$

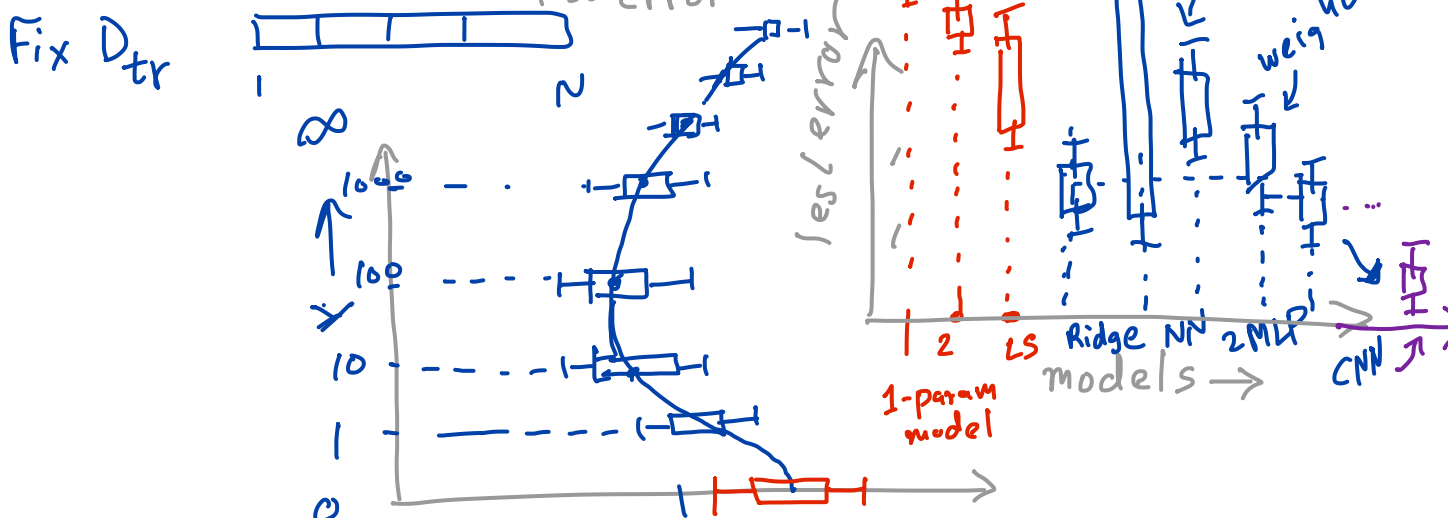
Ridge estimate

$\lambda = \infty$



Result of Cross-validation

bias decreases  
variance increases



$$y_n \approx \beta_0$$

$$y_n \approx \beta_0 + \beta_1 x_{n1}$$



# 10 Recent Advances

## Deep Learning & Overfitting

Deep learning has shown a new (but old) way to combat overfitting. For many applications, more data and deep architecture combined with stochastic gradient-descent is able to get us to a good minimum which generalizes well.

## Challenges

There are many challenges ahead. Learning from nasty, unreliable data still remains a challenge (e.g. small sample size, redundant data, non-stationary data, sequential learning).

On the other hand, living beings - even young ones - are very good in dealing with such data. How do they do it, and how can we design ML methods that can learn like them?