

Supervised learning  $\{x^{(i)} \in \mathbb{R}^D, t^{(i)} \in \mathbb{R}\}_{i=1}^N = \text{training set}$

→ linear regression

through ordinary least squares (OLS) / Residual Sum of Squares

general model  $h_{\beta}(x) = \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}$  loss

$$\text{loss } \mathcal{L}(h_{\beta}) = \mathcal{L}(\beta) = \frac{1}{2N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$

→ to learn  $h_{\beta}(x)$  from  $\{x^{(i)} \in \mathbb{R}^D, t^{(i)} \in \mathbb{R}\}_{i=1}^N$  (\*)

We can minimize (\*) through Gradient descent.

Concretely we start from random initial guess for  $\beta_j$

We then move iteratively in the direction  $d = -\text{grad}_{\beta} \ell(\beta)$

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \eta \text{grad}_{\beta} \ell(\beta)$$

$$\beta^{(k+1)} \leftarrow \beta^{(k)} - \eta \frac{1}{N} \sum_{i=1}^N (t^{(i)} - h_{\beta}(x^{(i)})) \cdot (-\tilde{x}^{(i)})$$

$$\tilde{x}^{(i)} = [1, x_1^{(i)}, \dots, x_D^{(i)}] \quad \beta = [\beta_0, \beta_1, \dots, \beta_D]$$

$$h_{\beta}(x) = \beta_0 + \sum_{j=1}^D \beta_j x_j = \beta^T \tilde{x} = [\beta_0 \ \beta_1 \ \dots \ \beta_D]$$

$$\begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_D \end{bmatrix}$$

2 main possibilities known as Batch gradient descent

(compute loss / gradient taking into account all the examples from the training set)

Stochastic gradient descent

(loop over the samples and perform one gradient step for each  $\{t^{(i)}, x^{(i)}\}$  pair)

in this setting, one pass over the whole training set corresponds to one epoch.

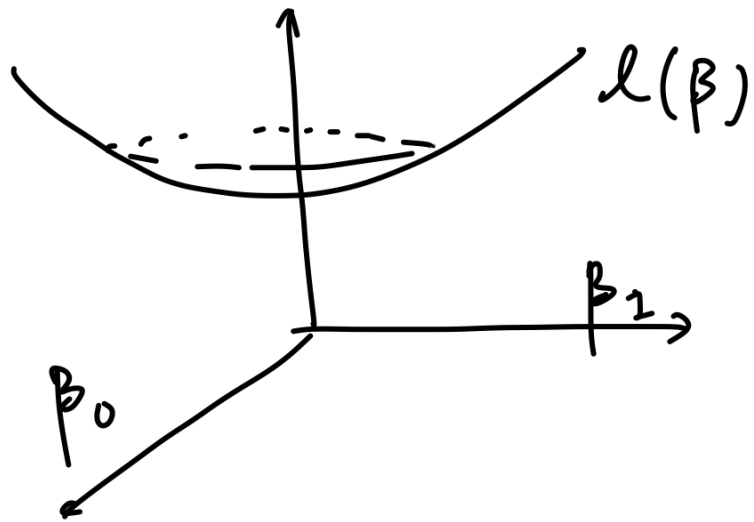
Today

- Alternative approach through Normal equations
- Statistical intuition on the least squares loss  
(Why is it a good idea to use the OLS loss?)
- Non linear data
- Bias Variance tradeoff
- Regularization

OLS loss: 
$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))^2$$
$$= \frac{1}{N} \sum_{i=1}^N e_i^2 = \frac{1}{N} \vec{e}^T \vec{e}$$

$$e_i = (t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)})), \quad \tilde{x}^{(i)} = [1, x^{(i)}]$$
$$= (t^{(i)} - \beta^T \tilde{x}^{(i)})$$

$$\vec{e} = [e_1, e_2, \dots, e_N]$$



## Approach #2

At the  $\beta$  that minimizes  $l(\beta)$   
the derivative must be zero

So in particular, to find the  $\beta$  that

minimizes  $l(\beta)$  we can compute the expression

of  $\text{grad } l(\beta) = \left[ \frac{\partial l}{\partial \beta_0}, \frac{\partial l}{\partial \beta_1}, \dots, \frac{\partial l}{\partial \beta_D} \right]$  and then

Solve the equation  $\text{grad } l(\beta) = 0$  for  $\beta$ .

Recall that

$$l(\beta) = e^T e = \sum_{i=1}^N (t^{(i)} - \beta^T \tilde{x}^{(i)})^2$$

$$e = \begin{bmatrix} e_1 \\ e_2 \\ \vdots \\ e_N \end{bmatrix} = \begin{bmatrix} t^{(1)} - \beta^T \tilde{x}^{(1)} \\ t^{(2)} - \beta^T \tilde{x}^{(2)} \\ \vdots \\ t^{(N)} - \beta^T \tilde{x}^{(N)} \end{bmatrix} = \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} - \begin{bmatrix} 1 & (x^{(1)})^T \\ 1 & (x^{(2)})^T \\ \vdots & \vdots \\ 1 & (x^{(N)})^T \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}$$

$$= \begin{bmatrix} t^{(1)} \\ \vdots \\ t^{(N)} \end{bmatrix} - \begin{bmatrix} 1 & x_1^{(1)} & x_2^{(1)} & \dots & x_D^{(1)} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_1^{(N)} & x_2^{(N)} & \dots & x_D^{(N)} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_D \end{bmatrix}$$
$$= \underline{\underline{T}} - \underline{\underline{X}} \underline{\underline{\beta}} \quad \rightarrow \quad \underline{\underline{X}}$$

$\vec{t}$  = target vector

$\tilde{X}$  = feature matrix (where rows are the feature vectors augmented with a '1')

From this,  $l(\beta)$  can read as

$$l(\beta) = \frac{1}{2N} e^T e = \frac{1}{2N} (\vec{t} - \tilde{X} \beta)^T (\vec{t} - \tilde{X} \beta)$$

$$\tilde{X} = \begin{bmatrix} \tilde{X}^{(1)} \\ \vdots \\ \tilde{X}^{(N)} \end{bmatrix}$$



$$L(\beta) = \frac{1}{2N} t^T t + \frac{1}{2N} \beta^T \underline{\tilde{X}}^T \underline{\tilde{X}} \beta - \frac{1}{2N} \underline{\tilde{t}}^T \underline{\tilde{X}} \beta - \frac{1}{2N} \beta^T \underline{\tilde{X}} \underline{\tilde{t}}$$

$$= \frac{1}{2N} \cancel{t^T t} + \frac{1}{2N} \beta^T \underline{\tilde{X}}^T \underline{\tilde{X}} \beta - \frac{1}{N} \underline{\tilde{t}}^T \underline{\tilde{X}} \beta$$

does not depend on  $\beta$  ②

$$\text{grad} = \left[ \frac{\partial L}{\partial \beta_0}, \dots, \frac{\partial L}{\partial \beta_D} \right]$$

term has the form

$$w^T \beta \text{ with}$$

$$w = \underline{\tilde{t}}^T \underline{\tilde{X}}$$

① gradient of  $w^T \beta = \underline{\tilde{t}}^T \underline{\tilde{X}} \beta$

$$= w_0 \beta_0 + w_1 \beta_1 + \dots + w_D \beta_D$$

$$\frac{\partial (w^T \beta)}{\partial \beta_0} = w_0 \quad \frac{\partial (w^T \beta)}{\partial \beta_1} = w_1 \quad \dots \quad \frac{\partial (w^T \beta)}{\partial \beta_D} = w_D$$

$$\rightarrow \text{grad}_{\beta} (\underline{\tilde{t}}^T \underline{\tilde{X}} \beta) = \left[ \frac{\partial \underline{\tilde{t}}^T \underline{\tilde{X}} \beta}{\partial \beta_0} \quad \dots \quad \frac{\partial \underline{\tilde{t}}^T \underline{\tilde{X}} \beta}{\partial \beta_D} \right]$$

$$= \underline{\tilde{X}}^T \underline{\tilde{t}} \leftarrow \textcircled{1}$$

term ② gradient of  $\beta^T \overbrace{\tilde{X}^T \tilde{X}}^{\text{Hessian}}$  with respect to  $\beta$  ?

$$\rightarrow \begin{bmatrix} \beta_0 & \beta_1 \end{bmatrix} \begin{bmatrix} H_{11} & H_{12} \\ H_{12} & H_{22} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} = \beta_0^2 H_{11} + \underbrace{\beta_1 \beta_0}_{\text{cross}} H_{12} + \underbrace{\beta_1^2}_{\text{cross}} H_{22} + \beta_1 \beta_0 H_{12}$$

$$\rightarrow \text{take } \frac{\partial \mathcal{L}}{\partial \beta_0} = 2\beta_0 H_{11} + 2\beta_1 H_{12}$$

$$\frac{\partial \mathcal{L}}{\partial \beta_1} = 2\beta_1 H_{22} + 2\beta_0 H_{12}$$

$$\text{grad } \mathcal{L} = \begin{bmatrix} \frac{\partial \mathcal{L}}{\partial \beta_0} \\ \frac{\partial \mathcal{L}}{\partial \beta_1} \end{bmatrix} = 2 \underbrace{H}_{\text{matrix}} \underbrace{\beta}_{\text{vector}}$$

For ② gradient is simply given by  $2 \underbrace{\tilde{X}^T \tilde{X}^T}_{\text{matrix}} \beta \frac{1}{2N}$

Combining ① & ② the gradient of  $\mathcal{L}(\beta)$  w.r.t  $\beta$  is given

$$\text{by } \frac{1}{N} \underbrace{\tilde{X}^T \tilde{X}}_{\text{matrix}} \beta - \frac{1}{N} \underbrace{\tilde{X}^T \vec{t}}_{\text{vector}} = \text{grad}_{\beta} \mathcal{L}(\beta)$$

In order to find the  $\beta$  that minimizes  $l(\beta)$  we set  $\text{grad } l(\beta)$  to 0 and solve the resulting equations

$$\frac{1}{N} \tilde{X}^T \tilde{X} \beta - \frac{1}{N} \tilde{X}^T \tilde{t} = 0$$

$$\tilde{X}^T \tilde{X} \beta = \tilde{X}^T \tilde{t}$$

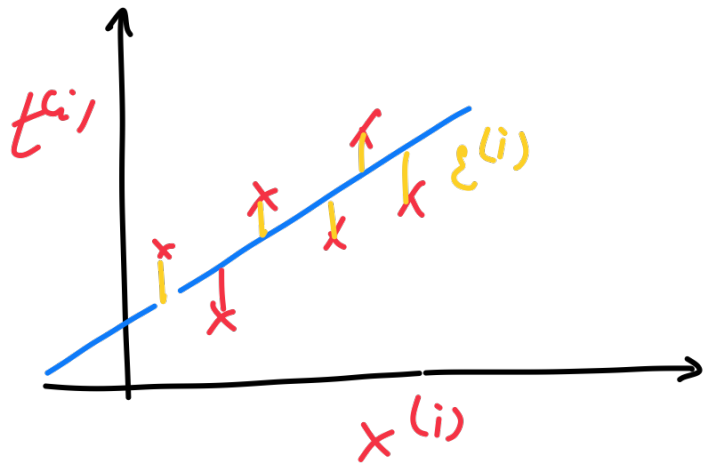
Normal equations

$$\beta = \left( \tilde{X}^T \tilde{X} \right)^{-1} \tilde{X}^T \tilde{t} \quad (**)$$

→ provided  $\tilde{X}^T \tilde{X}$  is invertible, we can compute  $\beta$  using (\*\*)

Why is it a good idea to choose  $\beta$  as the model that minimizes the least squares loss?

Assume ①  $t^{(i)} = \beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)} + \varepsilon^{(i)}$



Further assume ②  $\varepsilon^{(i)} \sim \mathcal{N}(0, \sigma)$

Finally assume ③  $\varepsilon^{(i)}$  independent.

(i.i.d  
independent  
identically distributed)

Only randomness in  $t^{(i)}$  is due to  $\varepsilon^{(i)}$ . In particular the probability of observing a particular  $t^{(i)}$  is given by

$$P(t^{(i)}) = P(\varepsilon^{(i)} = t^{(i)} - (\beta_0 + \beta_1 x_1^{(i)} + \dots + \beta_D x_D^{(i)}))$$

Assumptions  
①+②

$$= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\varepsilon^{(i)})^2}{2\sigma^2}\right) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t^{(i)} - \beta^T \tilde{x}^{(i)})^2}{2\sigma^2}\right)$$

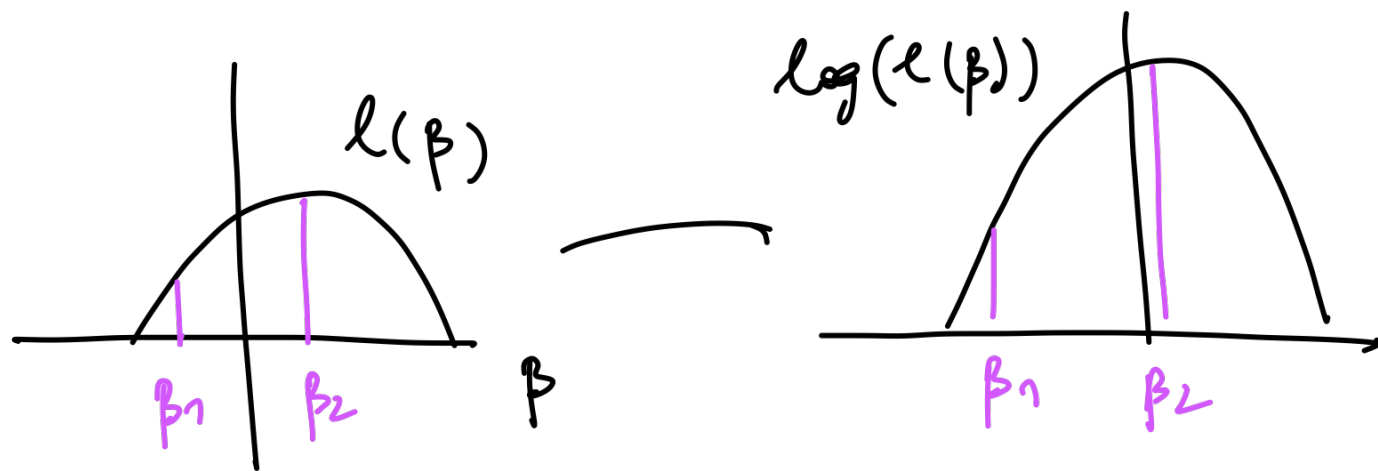
From Assumption ③ the probability of observing  $\{t^{(i)}\}_{i=1}^N$  given  $\{x^{(i)}\}_{i=1}^N$  is given by

When viewed as a function of  $\beta$   
= likelihood (\*)

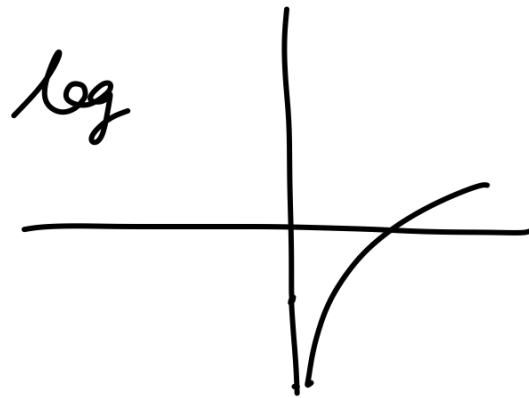
$$P(\{t^{(i)}\}_{i=1}^N | \{x^{(i)}\}_{i=1}^N) = \prod_{i=1}^N P(t^{(i)}) = \left( \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(t^{(i)} - \beta^T \tilde{x}^{(i)})^2}{2\sigma^2}\right) \right)^2$$

Given a set of observations  $\{t^{(i)}\}_{i=1}^N$  and a distribution on these observations  $p(\{t^{(i)}\}_{i=1}^N)$  that is parametrized by  $\beta$  a good approach is to look for the value of  $\beta$  that leads to  $p(\{t^{(i)}\}_{i=1}^N)$  being maximized.

in order to maximize (\*) with respect to  $\beta$  we first get rid of the product by taking the log



for  $\beta_1, \beta_2$   
 $l(\beta_1) < l(\beta_2)$   
 $\log(l(\beta_1)) <$   
 $\log(l(\beta_2))$



$\Rightarrow$  As  $\log$  is an increasing function, looking for  $\operatorname{argmax}_{\beta} l(\beta)$   
 is equivalent to  $\operatorname{argmax}_{\beta} \log(l(\beta))$



Applying this idea to (\*)

$$\log \left( \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} \exp \left( - \frac{(t^{(i)} - \beta^T \tilde{x}^{(i)})^2}{2\sigma^2} \right) \right) \rightarrow \text{log likelihood}$$

$$= \sum_{i=1}^N \log \left( \frac{1}{\sqrt{2\pi}\sigma} \right) - \frac{(t^{(i)} - \beta^T \tilde{x}^{(i)})^2}{2\sigma^2}$$

does not depend on  $\beta$

We can then find  $\beta^*$  by maximizing

$$LL = - \sum_{i=1}^N \frac{(t^{(i)} - \beta^T \tilde{x}^{(i)})^2}{2\sigma^2}$$

Maximum likelihood estimator for  $\beta$

$\hat{\beta}_{MLE}$

→ the  $\beta$  that maximize the likelihood is the  $\beta$  that minimizes the OLS loss

When the  $f_i$ 's are generated following assumptions ① to ③  
minimizing the OLS loss will return the Maximum Likelihood  
estimator for  $\beta$ .