



Lecture 7: Linear supervised learning



Nando de Freitas

www.cs.ubc.ca/~nando/340-2009/

September 2009

Outline

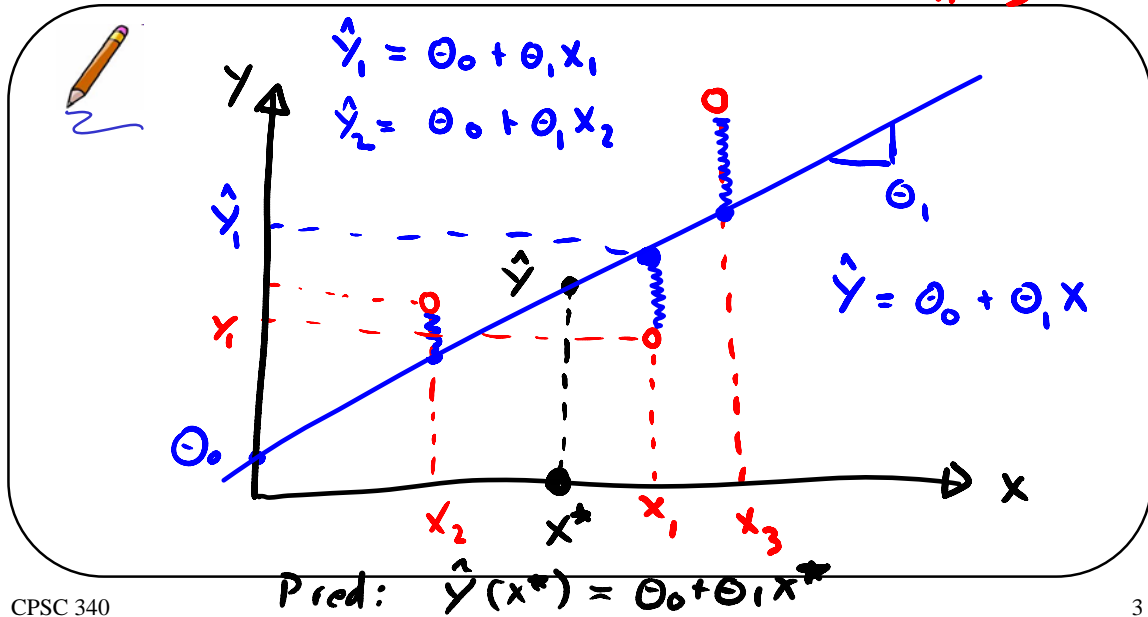
Linear regression is a supervised learning task. It is of great interest because:

- Many real processes can be approximated with linear models.
- Linear regression appears as part of larger problems.
- It can be solved analytically.
- It illustrates many of the approaches to machine learning.

Least squares

Given the data $\{x_{1:n}, y_{1:n}\}$, with $x_i \in \mathbb{R}^d$ and $y_i \in \mathbb{R}$, we want to fit a hyper-plane that maps x to y .

$d=1$
 $n=3$



CPSC 340

3


Least squares

$$\min \left\{ (y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 \right\}$$
$$= \min_{\theta_0, \theta_1} \left\{ (y_1 - \theta_0 - \theta_1 x_1)^2 + (y_2 - \theta_0 - \theta_1 x_2)^2 + (y_3 - \theta_0 - \theta_1 x_3)^2 \right\}$$

CPSC 340

4

Learning and prediction with least squares

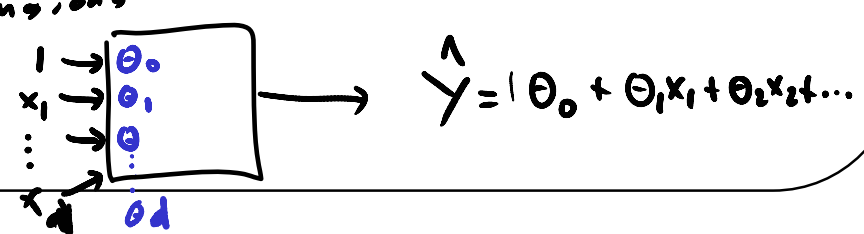
 Learning

Given $x_{1:n}, y_{1:n} \Rightarrow \hat{\Theta}$

prediction

$x^* \rightarrow \boxed{\hat{y} = \Theta x^*} = \hat{y}(x^*)$

in d dimensions


 $\hat{y} = 1\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$

CPSC 340 5

Least squares

Mathematically, the linear model is expressed as follows:

$i = 1, 2, \dots, n$

$\begin{pmatrix} x_{i,1}, y_i \\ \vdots \\ x_{i,n}, y_n \end{pmatrix}$

$\hat{y}_i = \theta_0 + \sum_{j=1}^d x_{ij} \theta_j$

eg. $i \equiv$ index over images
 $j \equiv$ index over features
 eg. (r, g, b)
 $d = 3$

We let $x_{i,0} = 1$ to obtain $\hat{y}_i = \sum_{j=0}^d x_{ij} \theta_j$

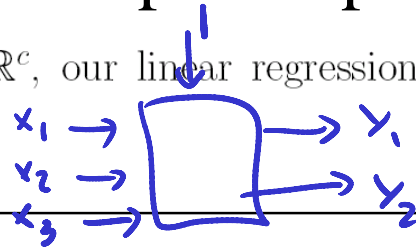
In matrix form, this expression is $\hat{Y} = X\theta$


$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} = \begin{bmatrix} x_{10} & \cdots & x_{1d} \\ \vdots & \vdots & \vdots \\ x_{n0} & \cdots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_0 \\ \vdots \\ \theta_d \end{bmatrix}$$

$y_1 = x_1 \theta_1 + \theta_0$
 $y_i = x_i \theta_1 + 1 \theta_0$
 $y_i = \begin{bmatrix} 1 & x_i \end{bmatrix} \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$

Least squares with multiple outputs

If we have several outputs $y_i \in \mathbb{R}^c$, our linear regression expression becomes:




 $Y = X\Theta$ $c = 2$ ~~$c = 3$~~
 $d = 3$

$$\begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ \vdots & \vdots \\ y_{n1} & y_{n2} \end{bmatrix} = \begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1d} \\ 1 & x_{21} & x_{22} & \dots & x_{2d} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nd} \end{bmatrix} \begin{bmatrix} \theta_{01} & \theta_{02} \\ \theta_{11} & \theta_{12} \\ \vdots & \vdots \\ \theta_{d1} & \theta_{d2} \end{bmatrix}$$

$n \times c$ $n \times (d+1)$ $(d+1) \times c$

Linear classification

 $\hat{Y} = X\Theta$ $n \times 3$ $n \times 3$ 3×3

data: $Y = \begin{bmatrix} 1 & x & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$

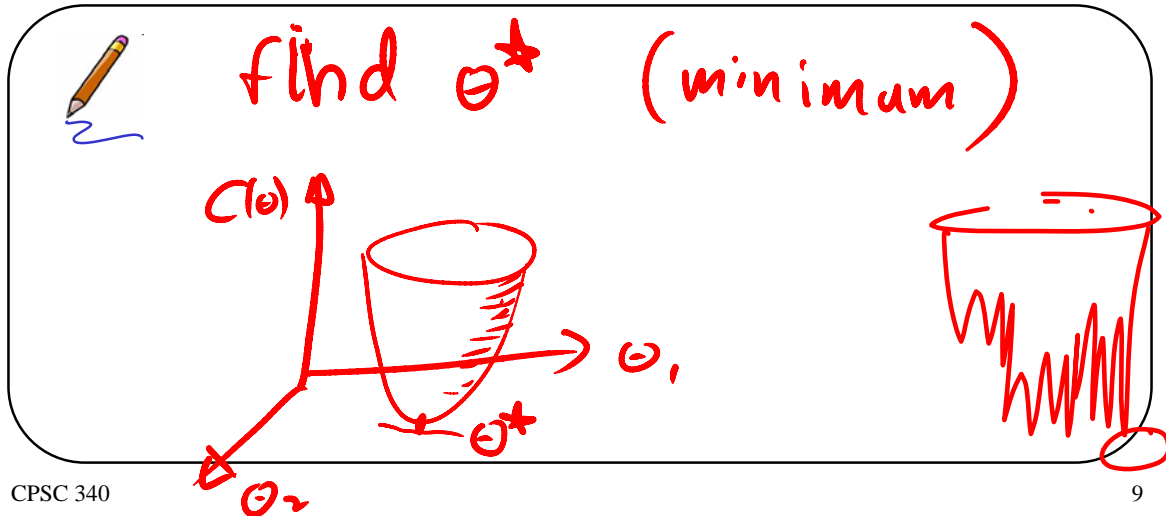
$\hat{Y} = X\Theta = \begin{bmatrix} 1 & 1.2 & 3.4 \\ 1 & 1.3 & 4 \\ \vdots & \vdots & \vdots \end{bmatrix} \begin{bmatrix} \theta_{10} & \theta_{11} & \theta_{12} \\ \theta_{20} & \theta_{21} & \theta_{22} \\ \theta_{30} & \theta_{31} & \theta_{32} \end{bmatrix}$

$\hat{y}(\bullet) = [0.8 \quad 0.2 \quad 0] \xrightarrow{\text{arg max}} [1 \quad 0 \quad 0]$

Optimization approach

Our aim is to minimise the quadratic cost between the output labels and the model predictions

$$C(\theta) = (Y - X\theta)^T(Y - X\theta)$$



Optimization approach

We will need the following results from matrix differentiation:
 $\frac{\partial A\theta}{\partial \theta} = A^T$ and $\frac{\partial \theta^T A \theta}{\partial \theta} = 2A^T \theta$ **EXAM**

$$\begin{aligned} \frac{\partial C}{\partial \theta} &= \frac{\partial}{\partial \theta} (Y - X\theta)^T (Y - X\theta) \\ &= \frac{\partial}{\partial \theta} (Y^T Y - \underbrace{2Y^T X \theta + \theta^T X^T Y}_{1 \times 1}) \\ &= 0 - 2X^T Y + 2X^T X \theta \rightarrow 0 \\ \hat{\theta} &= (X^T X)^{-1} X^T Y \end{aligned}$$

CPSC 340 10

Optimization approach



θ has $d+1$ parameters. For simplicity, I will say d .

These are the normal equations. The solution (estimate) is:

$$\hat{\theta} = \begin{matrix} d \times 1 & d \times d & d \times 1 \\ \underbrace{(X^T X)^{-1}} & \underbrace{X^T} & \underbrace{Y} \end{matrix}$$

The corresponding predictions are

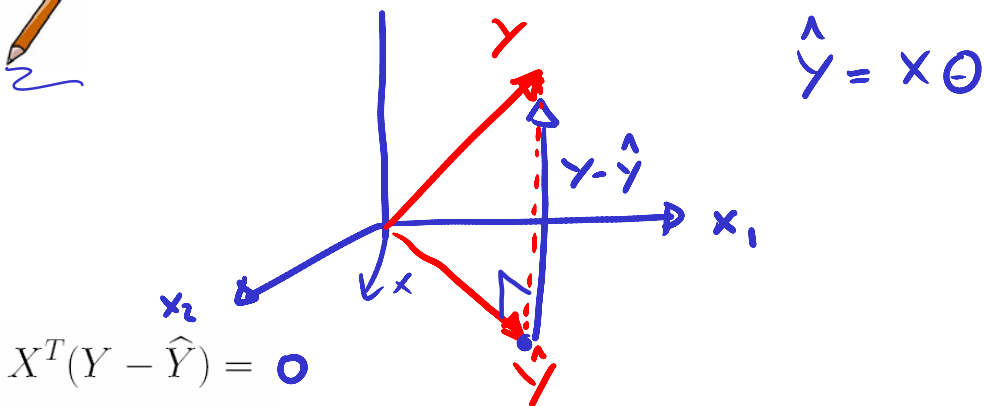
$$\hat{Y} = X \hat{\theta}$$

$$\hat{Y}_k = H Y = X_k (X^T X)^{-1} X^T Y$$

where H is the "hat" matrix. H

X_k is a new point for which I want a prediction.

Geometric approach



$$x^T y - x^T \hat{y} = 0$$

$$x^T y = x^T x \theta = 0$$

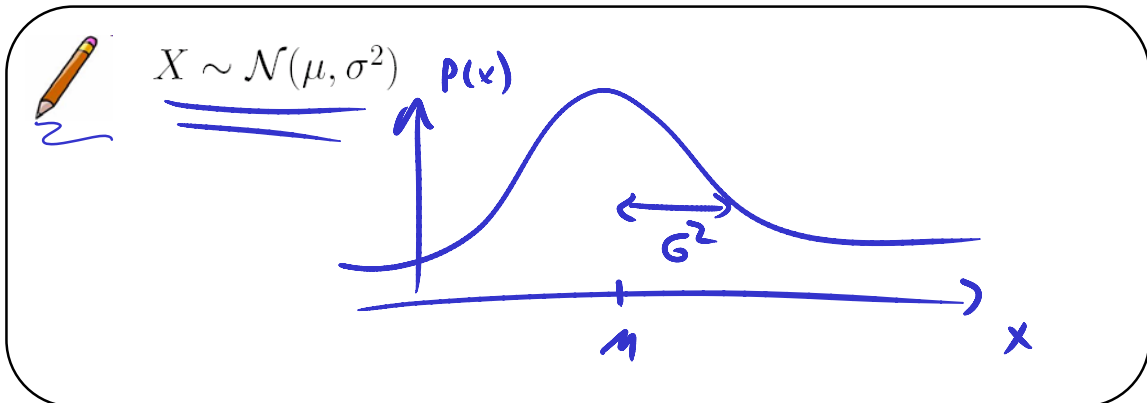
$$\theta = (x^T x)^{-1} x^T y$$

Probability approach: Univariate Gaussian distribution

The probability density function of a Gaussian distribution is given by

$$p(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2}$$

where μ is the mean or center of mass and σ^2 is the variance.



CPSC 340

13

Multivariate Gaussian distribution

Let $x \in \mathbb{R}^n$. The pdf of an n-dimensional Gaussian is given

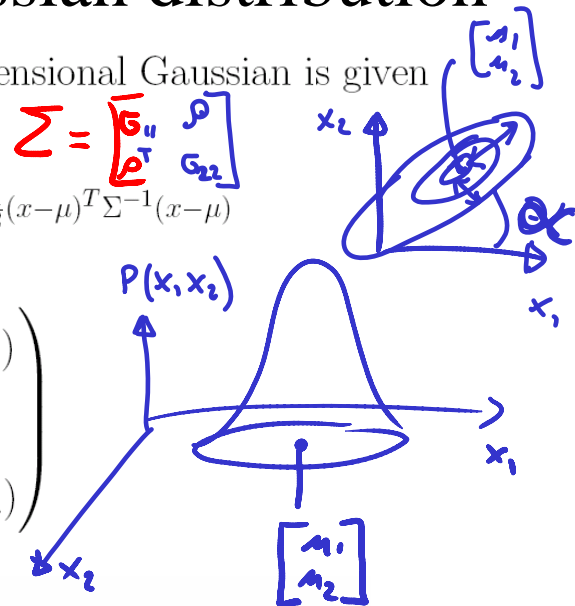
by

$$p(x) = \frac{1}{2\pi^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$

$$\mu = \begin{pmatrix} \mu_1 \\ \vdots \\ \mu_n \end{pmatrix} = \begin{pmatrix} \mathbb{E}(x_1) \\ \vdots \\ \mathbb{E}(x_n) \end{pmatrix}$$

$$\Sigma = \begin{pmatrix} \sigma_{11} & \cdots & \sigma_{1n} \\ \vdots & \ddots & \vdots \\ \sigma_{n1} & \cdots & \sigma_{nn} \end{pmatrix} = \mathbb{E}[(X - \mu)(X - \mu)^T]$$

$$\sigma_{ij} = \mathbb{E}[X_i - \mu_i)(X_j - \mu_j)^T]$$

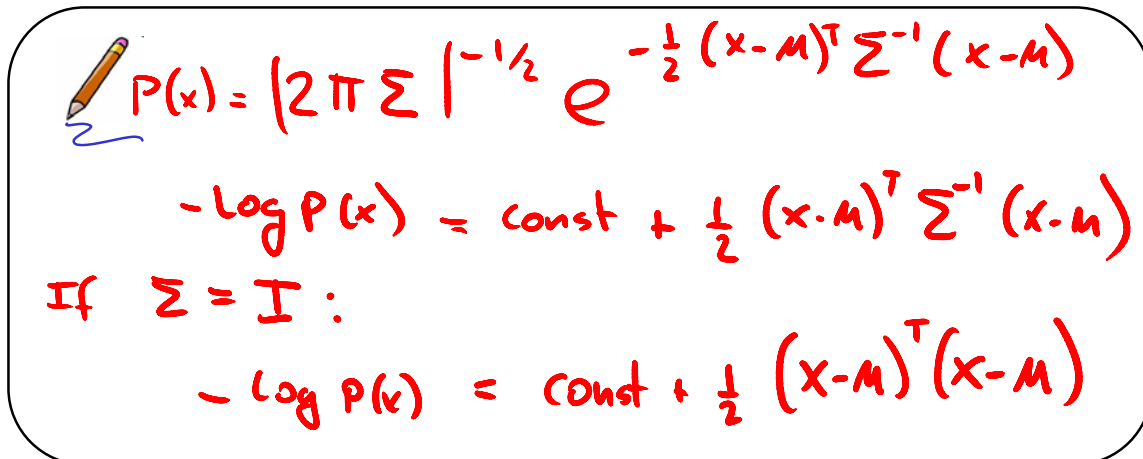


CPSC 340

14

Multivariate Gaussian distribution

We can interpret each component of x , for example, as a feature of an image such as colour or texture. The term $\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$ is called the **Mahalanobis distance**. Conceptually, it measures the distance between x and μ .


$$P(x) = (2\pi \Sigma)^{-1/2} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}$$
$$-\log P(x) = \text{const} + \frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)$$

If $\Sigma = I$:

$$-\log P(x) = \text{const} + \frac{1}{2}(x-\mu)^T (x-\mu)$$

CPSC 340

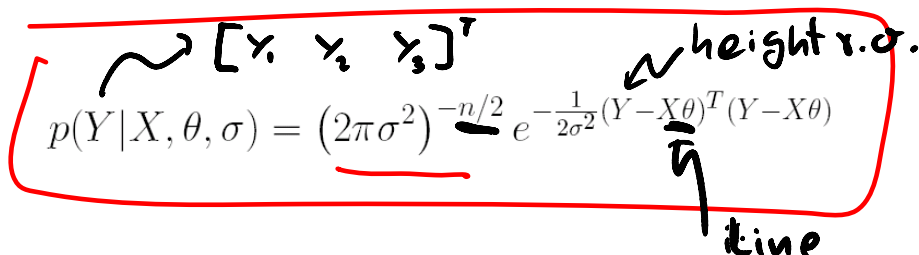
15

Maximum likelihood approach

If our errors are Gaussian distributed, we can use the model

$$Y = X\theta + \mathcal{N}(0, \sigma^2 I)$$
$$|\sigma^2 I_n| = (\sigma^2)^{n/2}$$

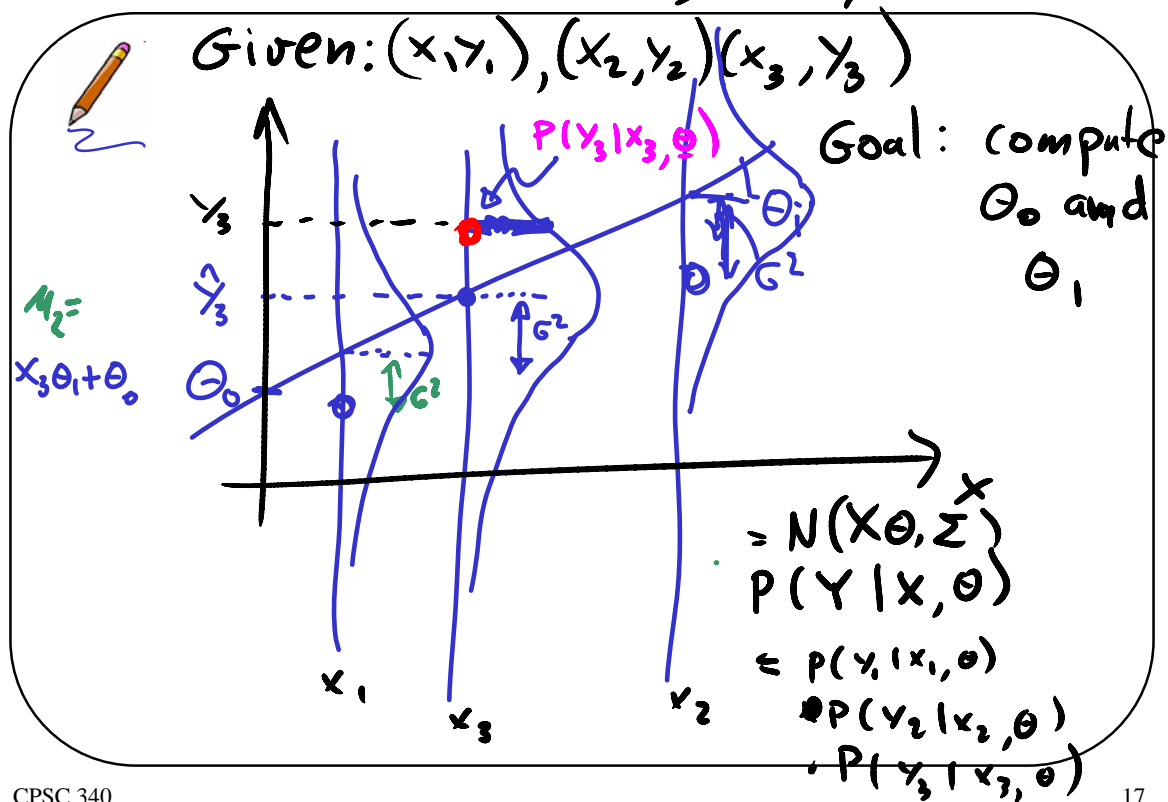
Note that the mean of Y is $X\theta$ and that its variance is $\sigma^2 I$. So we can equivalently write this expression using the probability density of Y given X , θ and σ : $\Sigma = \sigma^2 I$


$$p(Y|X, \theta, \sigma) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2}(Y-X\theta)^T (Y-X\theta)}$$

CPSC 340

16

Maximum likelihood ~~independent~~



Maximum likelihood

The maximum likelihood (ML) estimate of θ is obtained by taking the derivative of the log-likelihood, $\log p(Y|X, \theta, \sigma)$.

The idea of maximum likelihood learning is to maximise the likelihood of seeing some data Y by modifying the parameters (θ, σ) .

Maximum likelihood $-\frac{1}{2\sigma^2} (y-x\theta)^T (y-x\theta)$



The ML estimate of θ is:

$$P(y|x, \theta, \sigma^2) = (2\pi\sigma^2)^{-n/2} e^{-\frac{1}{2\sigma^2} (y-x\theta)^T (y-x\theta)}$$

$$\mathcal{L}(\theta) = \log P(y|x, \theta, \sigma^2) = \underbrace{-\frac{n}{2} \log(2\pi\sigma^2)}_{\text{Const. (ind. of } \theta)} - \frac{1}{2\sigma^2} (y-x\theta)^T (y-x\theta)$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta} = -\frac{1}{2\sigma^2} \frac{\partial}{\partial \theta} \left[(y-x\theta)^T (y-x\theta) \right] \rightarrow 0$$

$$\hat{\theta}_{LS} = \hat{\theta}_{ML} = (X^T X)^{-1} X^T Y$$

Maximum likelihood



Proceeding in the same way, the ML estimate of σ is:

$$\mathcal{L}(\sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (y-x\theta)^T (y-x\theta)$$

$$\frac{\partial \mathcal{L}(\sigma^2)}{\partial \sigma} = -\frac{n}{2} \frac{2\pi\sigma}{2\pi\sigma^2} + \frac{1}{2\sigma^3} (y-x\theta)^T (y-x\theta) \rightarrow 0$$

$$\frac{n}{\sigma} = \frac{1}{\sigma^3} (y-x\theta)^T (y-x\theta)$$

$$\sigma^2 = \frac{1}{n} (y-x\theta)^T (y-x\theta)$$

Lecture 8: Regularization and ridge regression



Nando de Freitas

www.cs.ubc.ca/~nando/340-2009/

September 2009

All the answers so far are of the form

$$\hat{\theta} = (X^T X)^{-1} X^T Y$$

They require the inversion of $X^T X$. This can lead to problems if the system of equations is poorly conditioned. A solution is to add a small element to the diagonal:

$$\hat{\theta} = (X^T X + \delta^2 I_d)^{-1} X^T Y$$

This is the ridge regression estimate. It is the solution to the following **regularised quadratic cost function**

$$C(\theta) = (Y - X\theta)^T (Y - X\theta) + \delta^2 \theta^T \theta$$

$$\textcircled{1} \frac{\partial Ax}{\partial x} = A^T$$

Proof

$$\frac{\partial x^T A x}{\partial x} = 2 A^T x$$



$$C(\theta) = (Y - X\theta)^T (Y - X\theta) + \delta^2 \theta^T \theta$$

$$\begin{aligned} C(\theta) &= Y^T Y - Y^T X \theta - \theta^T X^T Y + \theta^T X^T X \theta + \delta^2 \theta^T I \theta \\ &= Y^T Y - 2 Y^T X \theta + \theta^T (X^T X + \delta^2 I) \theta \end{aligned}$$

$$\begin{aligned} \frac{\partial C(\theta)}{\partial \theta} &= 0 - (2 Y^T X)^T + 2 (X^T X + \delta^2 I)^T \theta \\ &= -2 X^T Y + 2 (X^T X + \delta^2 I) \theta \end{aligned}$$

Equate to zero:

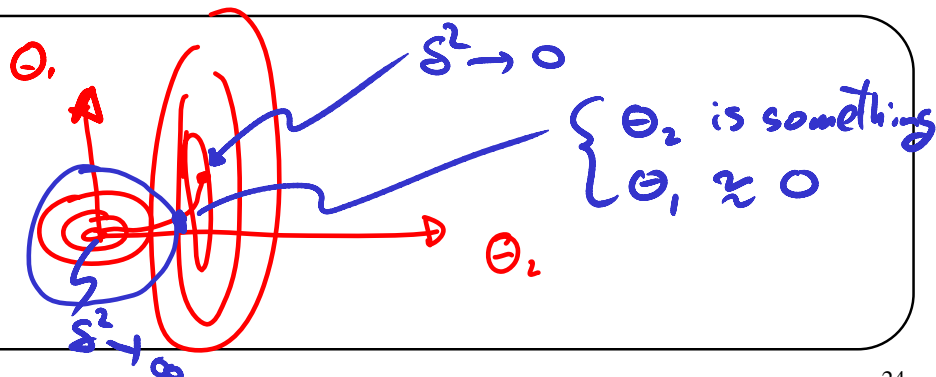
$$(X^T X + \delta^2 I) \theta = X^T Y$$

Ridge as constrained optimization

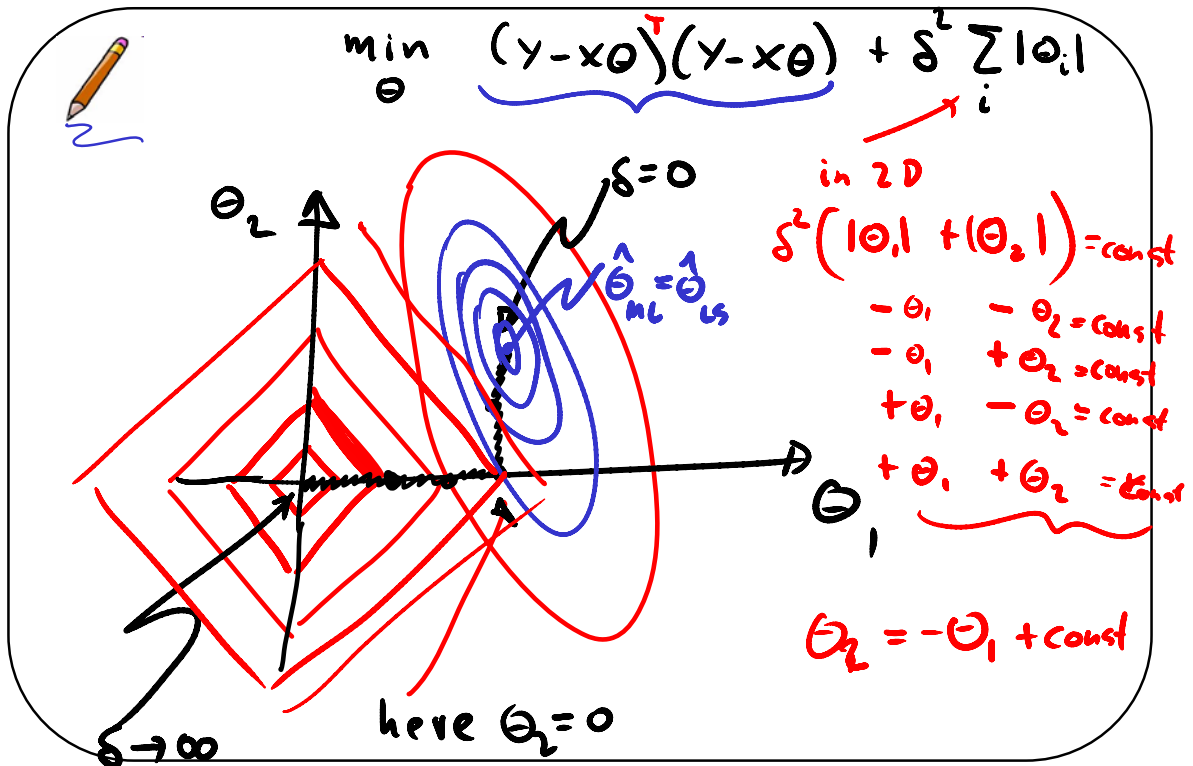
$$\min_{\theta : \theta^T \theta \leq t} \{(Y - X\theta)^T (Y - X\theta)\}$$

Large values of θ are penalised. We are **shrinking** θ towards zero. This can be used to carry out **feature weighting**.

An input $x_{i,d}$ weighted by a small θ_d will have less influence on the output y_i .



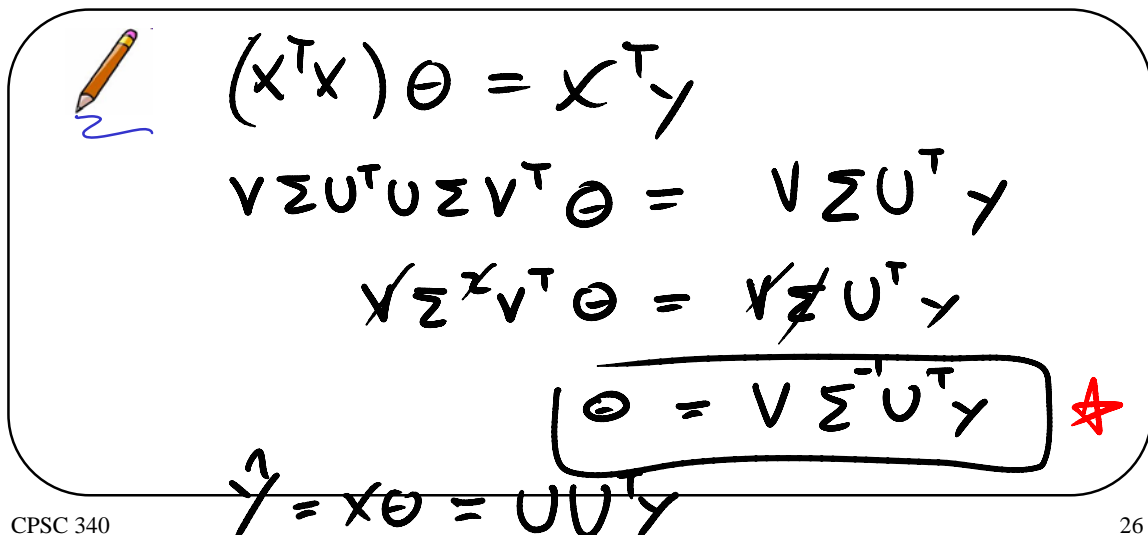
The Lasso



Spectral view of ridge regression

Again, let $X \in \mathbb{R}^{n \times d}$ be factored as $X = U\Sigma V^T = \sum_{i=1}^d u_i \sigma_i v_i^T$,

The least squares prediction is: $\hat{Y}_{LS} = \sum_{i=1}^d u_i u_i^T Y$





Likewise, for ridge regression we have:

$$\hat{Y}_{ridge} = \sum_{i=1}^d \frac{\sigma_i^2}{\sigma_i^2 + \delta^2} u_i u_i^T Y$$

in terms of SVD, what is
 $\hat{\Theta}_{ridge} = ?$

Regularization and noise filtering

The filter factor

$$f_i = \frac{\sigma_i^2}{\sigma_i^2 + \delta^2}$$

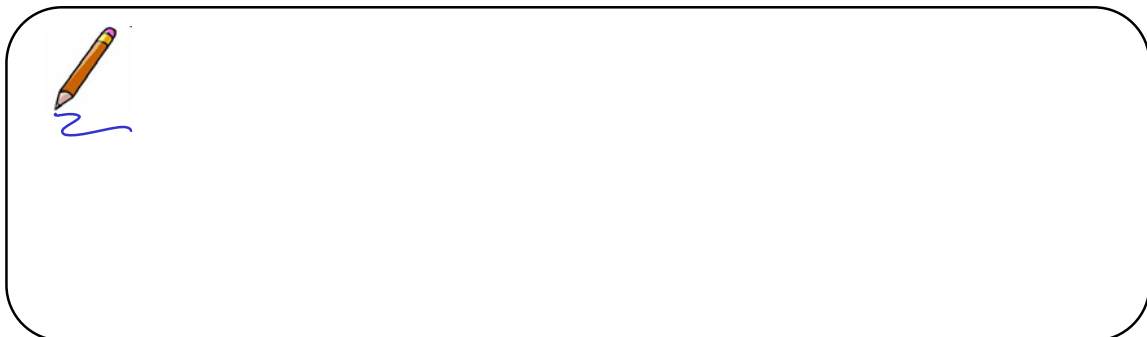
penalises small values of σ^2 (they go to zero at a faster rate).



Regularization and noise filtering

Small eigenvectors tend to be wobbly. The Ridge filter factor f_i gets rid of the wobbly eigenvectors. Therefore, the predictions tend to be more stable (smooth, regularised).

The smoothness parameter δ^2 is often estimated by cross-validation or Bayesian hierarchical methods.



$\hat{y} = X\theta$ Minimax and cross-validation

$\sum_{i \in \text{train}} (y_i - \hat{y}_i)^2$

δ^2	Error train	Error test	max	min max	avg.
0.1	100	2	100		
1	10	11	11	X	
10	1	19	19		X
50	20	0	20		X
100	100	1000	1000		

min max $\delta = 1$ // best avg $\delta = \begin{cases} 10 \\ 50 \end{cases}$

Lecture 9:

Bayesian learning for linear models



Nando de Freitas

www.cs.ubc.ca/~nando/340-2009/

September 2009

Bayesian linear-Gaussian supervised learning

In the Bayesian linear prediction setting, we focus on computing the posterior:

$$\begin{aligned} p(\theta|X, Y) &\propto p(Y|X, \theta)p(\theta) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2}(Y-X\theta)^T(Y-X\theta)} p(\theta) \end{aligned}$$

We often want to maximise the posterior — that is, we look for the *maximum a posteriori* (MAP) estimate. In this case, the choice of prior determines a type of constraint! For example, consider a Gaussian prior $\theta \sim \mathcal{N}(0, \delta^2\sigma^2 I_d)$. Then

$$p(\theta|X, Y) \propto (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2}(Y-X\theta)^T(Y-X\theta)} (2\pi\sigma^2\delta^2)^{-\frac{d}{2}} e^{-\frac{1}{2\delta^2\sigma^2}\theta^T\theta}$$



$$p(\theta|X, Y) = |2\pi\sigma^2 M|^{-\frac{1}{2}} e^{-\frac{1}{2\sigma^2}(\theta-\mu)^T M^{-1}(\theta-\mu)}$$
$$\propto (2\pi\sigma^2)^{-\frac{n}{2}} e^{-\frac{1}{2\sigma^2}(Y-X\theta)^T(Y-X\theta)} (2\pi\sigma^2\delta^2)^{-\frac{d}{2}} e^{-\frac{1}{2\delta^2\sigma^2}\theta^T\theta}$$



Bayesian posterior

So the posterior for θ is Gaussian:

$$p(\theta|X, Y) = |2\pi\sigma^2 M|^{-\frac{1}{2}} e^{-\frac{1}{2\sigma^2}(\theta-\mu)^T M^{-1}(\theta-\mu)}$$

with **sufficient statistics**:

$$\mathbb{E}(\theta|X, Y) = (X X^T + \delta^{-2} I_d)^{-1} X^T Y$$

$$\text{var}(\theta|X, Y) = (X X^T + \delta^{-2} I_d)^{-1} \sigma^2$$

Bayesian estimates, ridge and ML

The MAP point estimate is:

$$\hat{\theta}_{MAP} = (X X^T + \delta^{-2} I_d)^{-1} X^T Y$$

It is the same as the ridge estimate (except for a trivial negative sign in the exponent of δ), which results from the L_2 constraint. A flat (“vague”) prior with large variance (large δ) leads to the ML estimate.

$$\hat{\theta}_{MAP} = \hat{\theta}_{ridge} \xrightarrow{\delta^2 \rightarrow 0} \hat{\theta}_{ML} = \hat{\theta}_{SVD} = \hat{\theta}_{LS}$$