

5/18/23

Lecture 16

Multivariate Gaussian

Linear regression

D. D.

Posted

L VIII Linear
Regression

- Double Descent (soon)
 - for regression
- poll results !!

Multivariate Normal

$$S = \mathbb{R}^d \Rightarrow \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

Bivariate $d=2$

$$S = \mathbb{R}^2 \quad \begin{bmatrix} x \\ y \end{bmatrix}$$

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \Sigma\right) \Leftrightarrow x, y \text{ jointly Normal}$$

$$\mu \in \mathbb{R}^2$$

$$f_{xy}(x, y) = \frac{1}{2\pi|\Sigma|^{1/2}} e^{-\frac{1}{2} [x - \mu_x, y - \mu_y] \Sigma^{-1} \begin{bmatrix} x - \mu_x \\ y - \mu_y \end{bmatrix}}$$

Like $\sim \frac{1}{\sigma^2}$

$$v^\top \Sigma^{-1} v = \square \in \mathbb{R}$$

$$v^\top \Sigma^{-1} v \geq 0 \text{ when } \Sigma > 0$$

$$\text{Cov}(x, y) = E[(x - \mu_x)(y - \mu_y)] = 0 \text{ if } x, y \text{ independent}$$

Correlation coefficient

$$\rho_{xy} = \frac{\text{Cov}(x, y)}{\sqrt{\text{Var}(x)\text{Var}(y)}} \quad \begin{cases} > 0 \\ < 0 \end{cases}$$

var X

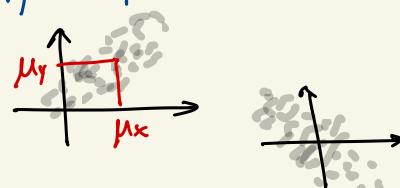
$\Sigma = \begin{pmatrix} \text{Var}(x) & \text{Cov}(x, y) \\ \text{Cov}(y, x) & \text{Var}(y) \end{pmatrix} \in \mathbb{R}^{2 \times 2}$

symmetric

$(\Sigma \succ 0) \quad \Sigma > 0$

positive definite

$\Leftrightarrow |\Sigma| > 0$



Multivariate $N(\mu, \Sigma)$

$$x = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$\mu = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix} \rightarrow \mu_i = E[X_i]$$

$$\Sigma = \begin{bmatrix} \sigma_1^2 & & & \\ \sigma_{1j} & \sigma_2^2 & & \\ & \ddots & \ddots & \\ & & & \sigma_d^2 \end{bmatrix}$$

with

$$\Sigma > 0 \iff v^T \Sigma v > 0 \text{ for any } v \neq 0$$

$$\sigma_{ij} = \sigma_{ji} = \text{Cov}(X_i, X_j)$$

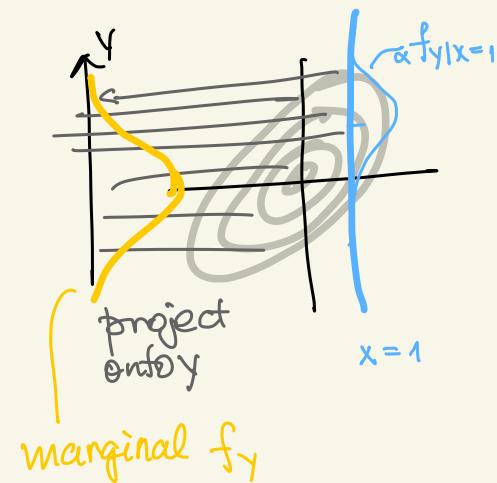
$$\sigma_i^2 = \text{Var } X_i$$

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

joint of (x_1, \dots, x_d)

$$(x-\mu) \Sigma^{-1} (x-\mu) = \square \geq 0$$

$$\sigma_{ij} = 0 \text{ for all } i \neq j : \Sigma = \begin{bmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \ddots & \ddots & 0 \\ \vdots & & \ddots & \sigma_d^2 \end{bmatrix} \Leftrightarrow x_1, x_2, \dots, x_d \text{ mutually independent}$$



Lecture Notes I – Regression

Marina Meilă
`mmp@stat.washington.edu`

Department of Statistics
University of Washington

May, 2023

Prediction problems



Linear regression



Linear regression for non-linear f

Reading: Ch.

Prediction

- ▶ **Data** $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots (x^n, y^n)\}$
- ▶ **Inputs** $x^{1:n} \in \mathbb{R}$, or \mathbb{R}^d
- ▶ **Outputs** $y^{1:n}$
- ▶ **Goal** Learn/estimate $f(x)$ **predictor** for y
- ▶ By type of output
 - ▶ **Classification** if $y \in S_Y$ discrete
 - ▶ $y \in \{0, 1\}$ (or $\{\pm 1\}$) **binary classification**
 - ▶ $y \in \{1, 2, \dots m\}$ **multiclass classification**
 - ▶ **Regression** if $y \in \mathbb{R}$ continuous

Prediction

- ▶ **Data** $\mathcal{D} = \{(x^1, y^1), (x^2, y^2), \dots, (x^n, y^n)\}$
- ▶ **Inputs** $x^{1:n} \in \mathbb{R}$, or \mathbb{R}^d
- ▶ **Outputs** $y^{1:n}$
- ▶ **Goal** Learn/estimate $f(x)$ **predictor** for y
- ▶ By type of output
 - ▶ **Classification** if $y \in S_Y$ discrete
 - ▶ $y \in \{0, 1\}$ (or $\{\pm 1\}$) **binary classification**
 - ▶ $y \in \{1, 2, \dots, m\}$ **multiclass classification**
 - ▶ **Regression** if $y \in \mathbb{R}$ continuous
- ▶ **Model family** $\mathcal{F} = \{f_\theta : \mathbb{R}^d \rightarrow S_Y, \theta \in \Theta\}$
 - ▶ $\mathcal{F} = \{ \text{linear functions} \}$ linear regression/classification
 - ▶ $\mathcal{F} = \{ \text{polynomes of degree } 2, 3, \dots \}$ polynomial regression/classification
 - ▶ $\mathcal{F} = \left\{ \frac{1}{1 + e^{-\beta^T x + \beta_0}} \right\}$ logistic regression
 - ▶ **neural network** regression/classification
 - ▶ **kernel** regression/classification
 - ▶ $\mathcal{F} = \{ \text{monotonic functions} \}$ **isotonic regression**
 - ▶ **support vector** regression/classification
 - ▶ regression/classification **trees** (and **random forests**)

Linear regression

→ predict y from x

Model is linear

$$\begin{cases} x \in \mathbb{R} \\ x \in \mathbb{R}^d \end{cases}$$

► $y^i = \beta_0 + \beta_1 x^i + \epsilon^i$ for $x^{1:n} \in \mathbb{R}$ (univariate regression)

► $y^i = \beta_0 + \beta_1 x_1^i + \beta_2 x_2^i + \dots + \beta_d x_d^i + \epsilon^i$ for $x^{1:n} \in \mathbb{R}^d$ (multivariate regression)

$y \in \mathbb{R}$ is output/response

$x_j^{1:n}$ for $j = 1 : d$ are input(s)/covariates/features/attributes/...

$\beta_{1:d}$ are regression coefficients, β_0 is intercept

$\epsilon^{1:n} \in \mathbb{R}$ is noise, $\epsilon^{1:n} \sim N(0, \sigma^2)$ i.i.d.

$$y^i = \beta_0 + \beta_1 x_1^i + \dots + \beta_d x_d^i + \varepsilon^i \quad \begin{matrix} \leftarrow \text{Model} \\ \text{noise} \end{matrix} \quad \varepsilon^i \sim N(0, \sigma^2) \quad \text{iid}$$

Data $\mathcal{D} \rightarrow \{(x^i, y^i)\}_{i=1:n}$

Wanted: parameters $\beta_0, \beta = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_d \end{bmatrix}, R^2$

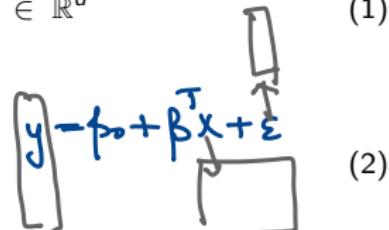
Linear regression model in matrix form

For a single data point (x^i, y^i)

$$\beta = \begin{bmatrix} \beta_1 \\ \beta_2 \\ \dots \\ \beta_d \end{bmatrix} \in \mathbb{R}^d \quad x^i = \begin{bmatrix} x_1^i \\ x_2^i \\ \dots \\ x_d^i \end{bmatrix} \in \mathbb{R}^d \quad (1)$$

Then,

$$y^i = \beta_0 + (x^i)^T \beta + \epsilon^i$$



For all \mathcal{D}

$$y = \begin{bmatrix} y^1 \\ y^2 \\ \dots \\ y^n \end{bmatrix} \in \mathbb{R}^n \quad X = \begin{bmatrix} (x^1)^T \\ (x^2)^T \\ \dots \\ (x^n)^T \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \epsilon = \begin{bmatrix} \epsilon^1 \\ \epsilon^2 \\ \dots \\ \epsilon^n \end{bmatrix} \quad (3)$$

$$y = \beta_0 \mathbf{1} + X\beta + \epsilon \quad \text{with } \text{Cov}(\epsilon) = \sigma^2 I_d \quad (4)$$



Solution by Maximum Likelihood

- ▶ **Likelihood** Probability $y | X$, parameters
- ▶ **Parameters** $\beta_0, \beta_{1:d}, \sigma^2$
- ▶ $\beta_0^{ML}, \beta_{1:d}^{ML}, (\sigma^2)^{ML} = \underset{\beta_0, \beta_{1:d}, \sigma^2}{\operatorname{argmax}} I(y|X, \beta_0, \beta_{1:d}, \sigma^2)$

The (log)-likelihood

- ▶ What is random? **the noise $\epsilon^{1:n}$**
- ▶ Express noise as function of $(x^{1:n}, y^{1:n})$

$$\text{for all } i: \underline{\epsilon^i} = y^i - \beta_0 - \beta^T x^i \sim N(0, \sigma^2) \quad (5)$$

- ▶ Likelihood

$$\text{Let } p_{0,\sigma^2}(\epsilon) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{\epsilon^2}{2\sigma^2}} = N(\epsilon; 0, \sigma^2)$$

- ▶ Then

$$L(\beta_0, \beta_{1:d}, \sigma^2) = \prod_{i=1}^n p_{0,\sigma^2}(\epsilon^i) \quad (6)$$

$$= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(\epsilon^i)^2}{2\sigma^2}} \quad \text{with } \epsilon^i$$

$$= \prod_{i=1}^n \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(y^i - \beta_0 - \beta^T x^i)^2}{2\sigma^2}} \quad (8)$$

- ▶ **log-likelihood**

$$l(\beta_0, \beta_{1:d}, \sigma^2) = \quad (9)$$

$$= \sum_{i=1}^n \left\{ -\frac{1}{2} \ln \sigma^2 - \frac{1}{2} \ln(2\pi) - \frac{1}{2} (y^i - \beta_0 - \beta^T x^i)^2 \frac{1}{2\sigma^2} \right\} \quad (10)$$

$$= -\frac{n}{2} \ln \sigma^2 - \frac{1}{2} \ln(2\pi) + \text{constant} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2 \quad (11)$$

Maximizing the log-likelihood w.r.t β

Calculus

$$\hat{\beta}^+ \text{ d } \boxed{w}$$

$$\hat{\beta}^T \hat{\beta} = I_d$$

~~$$\hat{\beta} \hat{\beta}^T = I_n$$~~

(12)

- For simplicity, let $\beta_0 = 0$; hence $y^i = \beta^T x^i + \epsilon^i$
- log-likelihood

$$I(\beta_{1:d}, \sigma^2) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^i - \beta^T x^i)^2 + \text{constant}$$

- For any σ^2 ,

$$\underset{\beta}{\operatorname{argmax}} I(\sigma^2, \beta) = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^n (y^i - \beta_0 - \beta^T x^i)^2$$

a Least Squares Problem

- In matrix form $\min_{\beta} \|y - X\beta\|^2$

- Solution

$$= (y - X\beta)^T (y - X\beta)$$

$$\beta^{ML} = (X^T X)^{-1} X^T y \quad X^+ \text{ pseudo inverse}$$

- with $(X^T X)^{-1} X^T \equiv X^+$ the **pseudoinverse** of X
- β^{ML} is **linear** in y !

(14)

\approx solve $X\beta = y$
with error

$$\begin{matrix} n & | & X \\ & | & \beta \\ d & | & y \end{matrix} \Rightarrow \begin{matrix} \text{solution} \\ \beta = X^+ y \end{matrix}$$

Statistical properties of β^{ML}

- ▶ Assume the true model is $y^i = \beta^T x^i + \epsilon^i$ with $\epsilon \sim N(0, \sigma^2)$.
- ▶ Here β, σ^2 are **true parameters** and we assume we know them.
- ▶
- ▶ β^{ML} is a random variable. Let us calculate its mean and standard deviation.
- ▶ **Expectation** of β^{ML}

$$E[\beta^{ML}] = E[X^\dagger y] = E[X^\dagger(X\beta + \epsilon)] \quad (15)$$

$$= E[X^\dagger X]\beta + E[X^\dagger \epsilon] \quad (16)$$

$$= \underbrace{X^\dagger X}_{I_d} \beta + X^\dagger \underbrace{E[\epsilon]}_0 = \beta \quad (17)$$

Hence, $E[\beta^{ML}] = \beta$ and we say that β^{ML} is **unbiased**

- ▶ **Covariance** of β^{ML}
- ▶ By a similar (but longer) calculation, we obtain that

$$\text{Cov}(\beta^{ML}) = \sigma^2(XX^T)^{-1} \in \mathbb{R}^{d \times d} \quad (18)$$

- ▶ In fact, β^{ML} has a Normal distribution. Remember from (15)

$$\beta^{ML} = \beta + X^\dagger \epsilon. \quad (19)$$

This is a linear transformation of ϵ , a Gaussian variable. Therefore,

$$\beta^{ML} \sim N(\beta, \sigma^2(XX^T)). \quad (20)$$

- ▶ This is a **multivariate Normal** distribution over \mathbb{R}^d ,

Estimation of σ^2

- Let $\hat{\epsilon}^i = y^i - (x^i)^T \beta^{ML}$, for $i = 1 : n$
- $\hat{\epsilon}^i$ called the residuals
- If we plug in β^{ML} in the log-likelihood equation (12) we obtain

$$I(\sigma^2, \beta^{ML}) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n (\hat{\epsilon}^i)^2 + \text{constant}$$

like $(x - \hat{\mu})^2$ for $N(\mu, \sigma^2)$ (21)

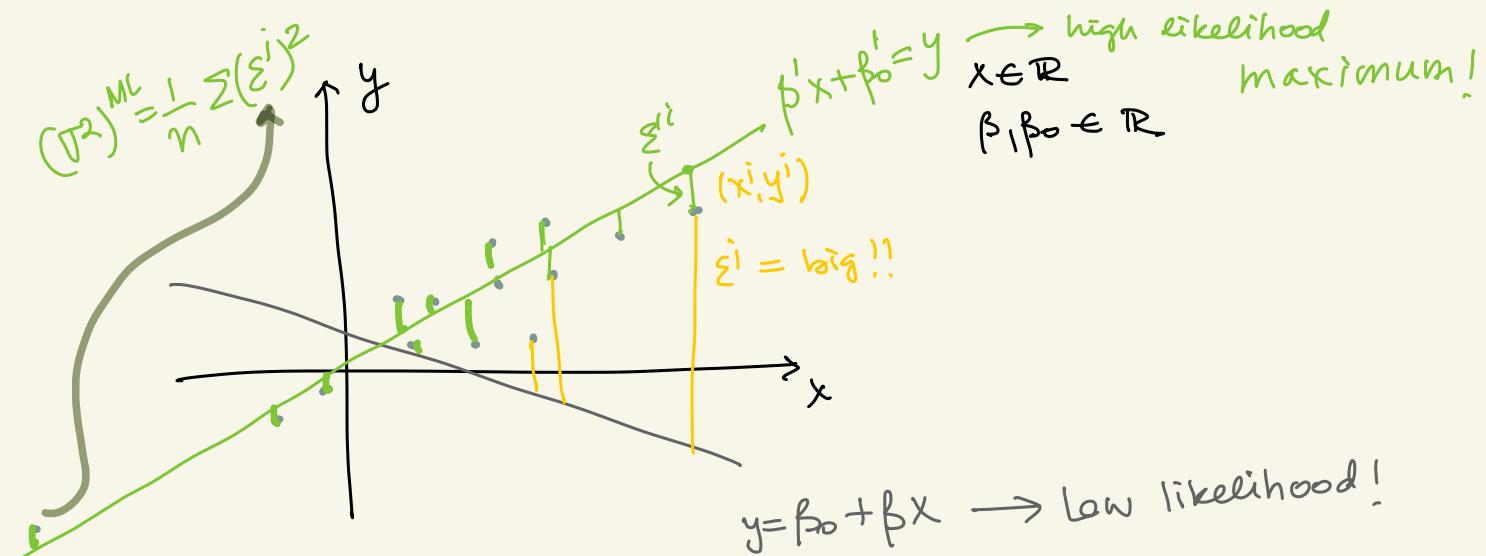
\uparrow known \uparrow known

- Maximizing this expression w.r.t. σ^2 gives

$$(\sigma^2)^{ML} = \frac{1}{n} \sum_{i=1}^n (\hat{\epsilon}^i)^2 \quad (22)$$

- The predictor is $f(x) = x^T \beta^{ML}$
- Hence, for a new $x \in \mathbb{R}^d$, our guess of y is $\hat{y} = f(x)$

new x : $\hat{f}(x) = (\beta^{ML})^T x$ \leftarrow guess for new y !

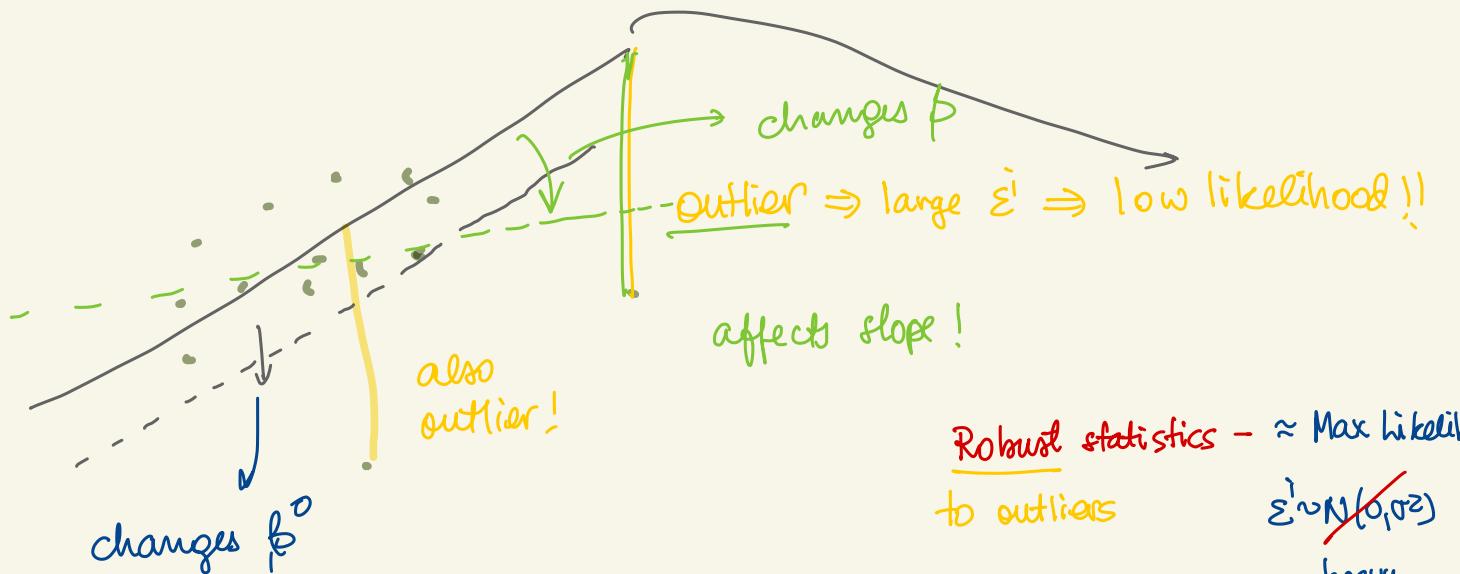


$$\begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix} = \begin{bmatrix} \beta_0 \\ \vdots \\ \beta_0 \end{bmatrix} + \begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix} \begin{bmatrix} \beta_1 \end{bmatrix} + \begin{bmatrix} \varepsilon^1 \\ \vdots \\ \varepsilon^n \end{bmatrix}$$

$$\begin{aligned}
 \beta^{ML} &= X^+ y = \frac{1}{\sum x_i^2} \begin{bmatrix} x^1 & \dots & x^n \end{bmatrix} \begin{bmatrix} y \\ \vdots \\ y^n \end{bmatrix} \\
 &= \frac{1}{(\sum x_i^2)} \sum_{i=1}^n x_i y_i
 \end{aligned}$$

if $\beta_0 = 0$:

$$X^+ = (X^T X)^{-1} X = \left(\sum_{i=1}^n (x^i)^2 \right)^{-1} \begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix} = \frac{1}{\sum_{i=1}^n (x^i)^2} \begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix}$$



Robust statistics - \approx Max likelihood
to outliers

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

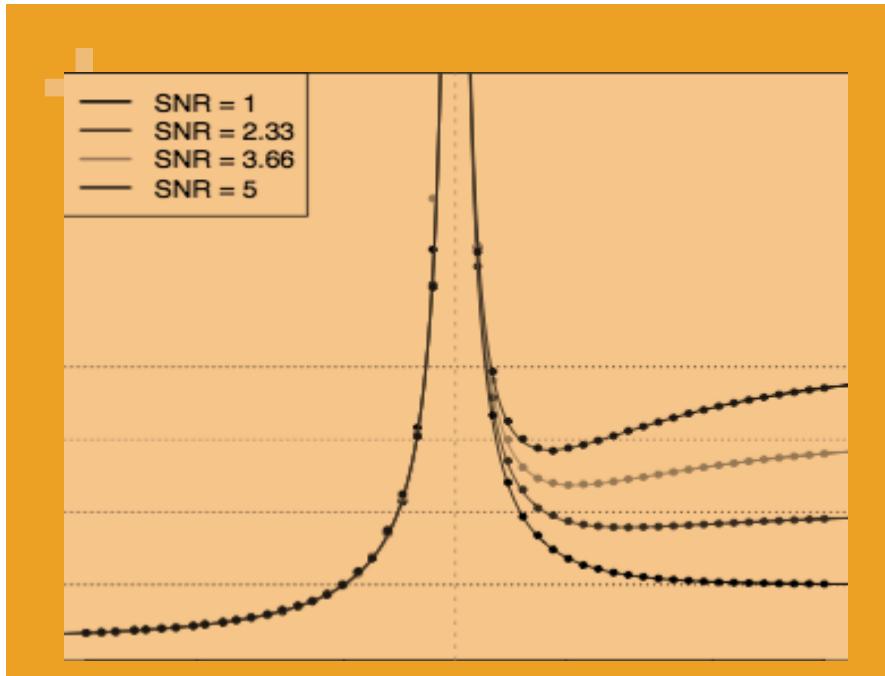
heavy tailed model

Linear regression → Linear regr
for non-linear
models

Logistic regression

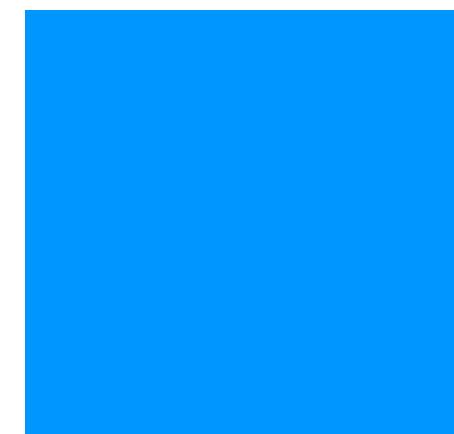
Overparameterized
(linear) regression

Bias - Variance
Tradeoff
for Neural
Network



Double Descent

Beyond the Bias-Variance trade-off



STAT 535+LPL2019

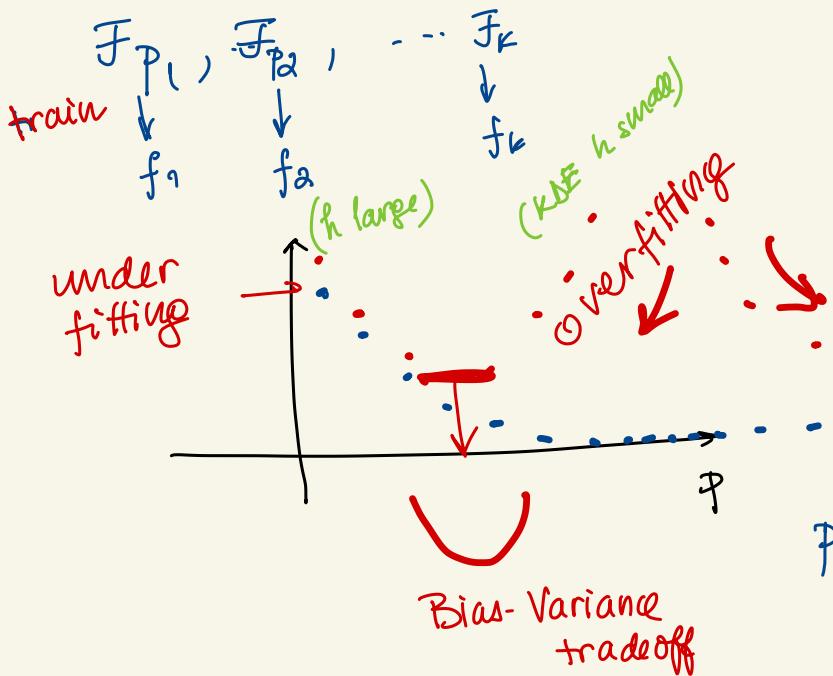
Marina Meila
University of Washington

Prediction with "overparametrized" model

$$S = \mathbb{R}^d \ni x$$

$$\mathcal{F}_P = \{ \text{neural network, RForest, } \dots \}_{f_\theta}$$

$$|\theta| = p \quad \# \text{ parameters}$$



$$\mathcal{D} = \{(x^i, y^i)\}_{i=1:n}$$

$\mathcal{D}_{\text{test}}$ large

trained on \mathcal{D}

$$\text{Least squares}$$

$$\text{Training error} = \sum_{i=1}^n (y^i - f(x^i))^2$$

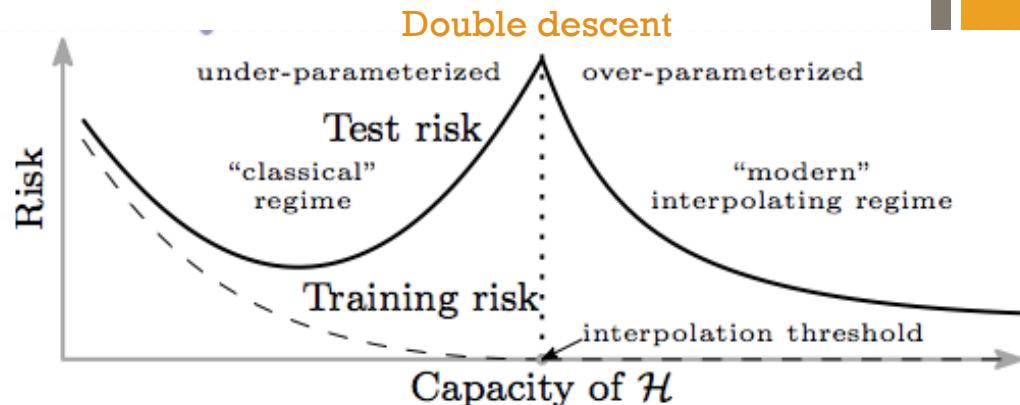
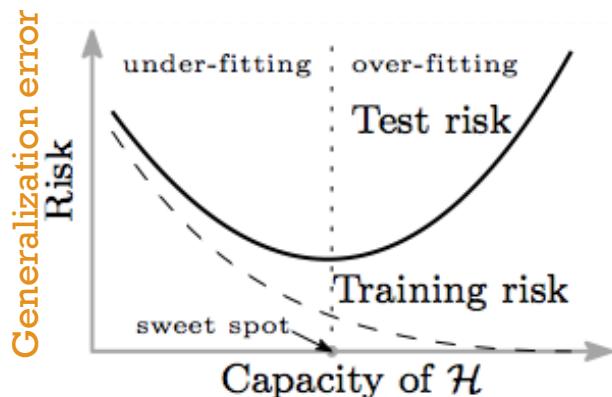
$$\text{Testing err} = \sum_{i \in \mathcal{D}_{\text{test}}} (y^i - f(x^i))^2$$

Double descent

$$p_1 < p_2 < \dots < p_K$$



What is observed



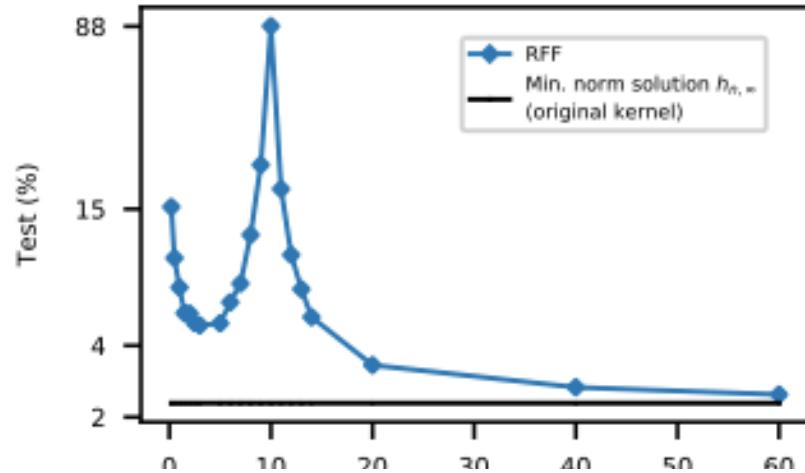
Belkin, Hsu, Ma, Mandal 2018

- Classical regime $p < N$
- Modern/Deep Learning/High dimensional regime $N > n$
 - Think N fixed, p increases, $\gamma = p/N$
 - Training error = 0 (interpolation)
 - **Test error decreases** with p (or γ)

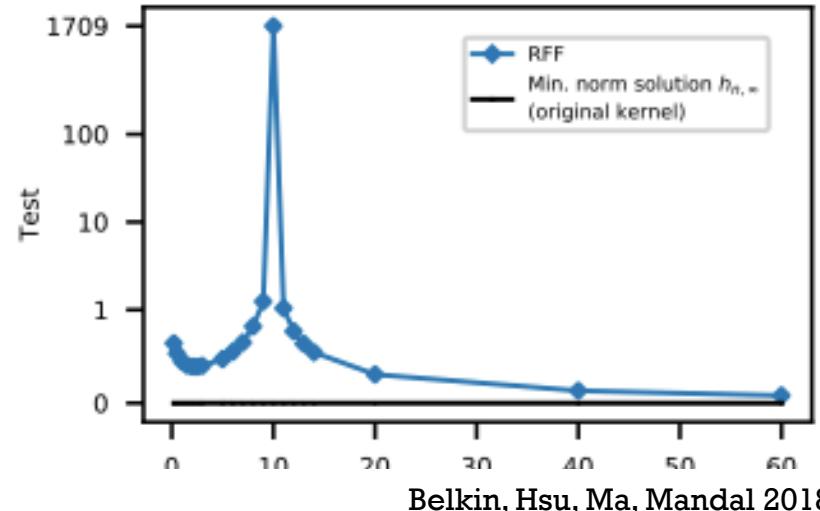


What is observed

Zero-one loss



Squared loss



- Double descent curves for the generalization error
 - Random Fourier Features (RFF)
 - ReLU 2 layer networks (with random first layer weights)
 - Random Forests, l2-Adaboost
 - Linear regression
- With and without noise

$$\begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} \beta_0 \end{bmatrix} + \begin{bmatrix} x^1 \\ \vdots \\ x^n \end{bmatrix} \begin{bmatrix} \beta \end{bmatrix} + \begin{bmatrix} \varepsilon^1 \\ \vdots \\ \varepsilon^n \end{bmatrix} \Rightarrow \begin{bmatrix} \varepsilon^1 \\ \vdots \\ \varepsilon^n \end{bmatrix} = \begin{bmatrix} y^1 \\ \vdots \\ y^n \end{bmatrix} - X\beta$$

$\underbrace{\quad}_{\text{dummy variable}}$

$$\begin{bmatrix} x^1 & 1 \\ x^2 & 1 \\ \dots & \\ x^n & 1 \end{bmatrix} \underbrace{\begin{bmatrix} \beta \\ \beta_0 \end{bmatrix}}_{\beta}$$

$$\boxed{\begin{array}{cccc} x^1 & \dots & x^n \\ \downarrow & \dots & \downarrow \end{array}}$$

$$X^+ = (X^T X)^{-1} X^T$$

$$X^T X = \begin{bmatrix} \sum_{i=1}^n (x^i)^2 & \sum_{i=1}^n x^i \\ \sum_{i=1}^n x^i & n \end{bmatrix} \in \mathbb{R}^{2 \times 2} \geq 0$$

invertible if X is full rank