

2/11/15

①

GAUSSIANS & REGRESSION

① High-level Picture & Flow of topic in Class

Last Time

Today / Next lecture

Real-World Problem: Binary / Discrete decisions
Nadal W/L, Email Spam/Non-spam

Continuous-valued variables
House prices, stock index, etc

Variable: $Y = y \in \{0, 1\}$ [or $\{1, \dots, k\}$]

$Y = y \in \mathbb{R}$ [later \mathbb{R}^d]

Model: $Y \sim \text{Ber}(\theta)$ [or $Y \sim \text{Cat}(\vec{\theta})$]

$Y \sim N(\mu, \sigma^2)$

Parameters: $\theta = P(Y=1) = E[\mathbb{I}\{Y=1\}]$

$\mu = E[Y]$ $\sigma^2 = \text{Var}(Y) = E[(Y-\mu)^2]$

MLE: $\hat{\theta}_{MLE} = \frac{\#H}{\#H + \#T}$

$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n y_i$ $\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$

Prior: $\theta \sim \text{Beta}(\beta_H, \beta_T) \Leftrightarrow P(\theta | \beta_H, \beta_T) \propto \theta^{\beta_H-1} (1-\theta)^{\beta_T-1}$

$\mu \sim N(\mu, \sigma^2)$ $\sigma \sim \text{Inverse Gamma}$

These are all models for Y (not too interesting)

We will later move to models of $Y | X=x$

output ↑ input/features

$Y | X=x \sim \text{Ber}(\theta_x)$
depends on input

$Y | X=x \sim N(\mu_x, \sigma_x^2)$

Classification

{ Logistic Regression
Naive Bayes }

Regression

{ Least Squares Regression }

② Meet the Gaussian Density Function

$$Y \sim N(\mu, \sigma^2)$$

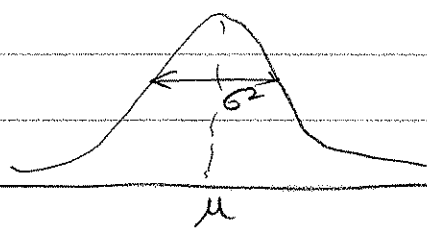
$$\equiv p(y | \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(y-\mu)^2}{2\sigma^2}}$$

Recall: Density function: $\int_{-\infty}^{\infty} p(y) dy = 1$

$$P_a[a \leq Y \leq b] = \int_a^b p(y) dy$$

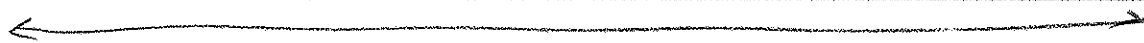
$$E[f(Y)] = \int_{-\infty}^{\infty} f(y) p(y) dy$$

$p(y) \geq 0$ but NOT necessarily < 1



$$\mu = E[Y]$$

$$\sigma^2 = E[(Y - \mu)^2]$$



③ Max-likelihood Estimation of μ, σ
(Same Principles)

may be exam scores of students

$$\text{Dataset } D = \{y_1, \dots, y_N\} \quad Y \sim N(\mu, \sigma^2)$$

$$(\hat{\mu}_{MLE}, \hat{\sigma}_{MLE}) = \underset{\mu, \sigma}{\operatorname{argmax}} LL(\mu, \sigma)$$

(2)

[Cooking through the math]

$$\begin{aligned}
 LL(\mu, \sigma) &= \log P(D | \mu, \sigma) \\
 &= \log \prod_{i=1}^N p(y_i | \mu, \sigma) \quad [\text{IID}] \\
 &= \sum_{i=1}^N \log p(y_i | \mu, \sigma) \\
 &= \sum_{i=1}^N \log \left[\frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - \mu)^2}{2\sigma^2}} \right] \\
 &= \sum_{i=1}^N \left[-\frac{(y_i - \mu)^2}{2\sigma^2} - \frac{1}{2} \log(2\pi\sigma^2) \right]
 \end{aligned}$$

Now, calculate!

$$\frac{\partial LL}{\partial \mu} = -\sum_{i=1}^N \frac{2(y_i - \mu) \cdot (-1)}{2\sigma^2} = 0$$

$$\Rightarrow \sum_{i=1}^N (y_i - \mu) = 0 \Rightarrow \boxed{\mu_{MLE} = \frac{1}{N} \sum y_i}$$

$$\frac{\partial LL}{\partial \sigma} = -\sum \frac{(y_i - \mu)^2}{2} \left(\frac{-2}{\sigma^3} \right) - \frac{N}{2} \frac{1}{2\pi\sigma^2} \cdot 2\pi \cdot (2\sigma) = 0$$

$$= \frac{\sum_{i=1}^N (y_i - \mu)^2}{\sigma^3} - \frac{N}{\sigma} = 0$$

$$\Rightarrow \frac{1}{\sigma_{MLE}^2} = \frac{1}{N} \sum_{i=1}^N (y_i - \mu_{MLE})^2 \quad \left\{ \text{assuming } \sigma \neq 0 \right\}$$

(4) MAP estimation of μ

We said $\{y_1, \dots, y_n\}$ may be exam scores of students (say)

$P(\mu)$ = I know I set the exam hard.
Class average should be around 20/100
(Just a thought. Don't panic)

For simplicity, assume $\sigma =$ (unknown) constant

$$\mu \sim N(\eta, \lambda^2)$$

↑ ↑
hyper-parameters
[known!]

Why another Gaussian?
It's self-conjugate!
(aka math is simple)

$$\hat{\mu}_{\text{MAP}} = \underset{\mu}{\text{argmax}} \log P(\mu | D)$$

$$= \underset{\mu}{\text{argmax}} \log \frac{P(D | \mu) P(\mu | \eta, \lambda^2)}{P(D)}$$

$$= \underset{\mu}{\text{argmax}} \log P(D | \mu) + \log P(\mu | \eta, \lambda^2) - \underbrace{\log P(D)}_{\text{const w.r.t } \mu}$$

$$\frac{\partial [\log P(\mu | D)]}{\partial \mu} = \frac{\partial \log P(D | \mu)}{\partial \mu} + \frac{\partial \log P(\mu | \eta, \lambda^2)}{\partial \mu}$$

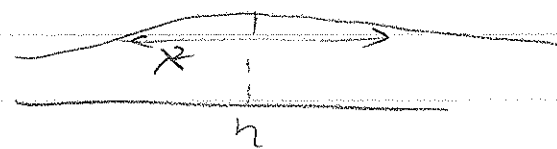
$$= \sum_{i=1}^N \frac{(y_i - \mu)}{\sigma^2} + \frac{\partial}{\partial \mu} \left[\frac{-(\mu - \eta)^2}{2\lambda^2} - \frac{1}{2} \log(2\pi\lambda^2) \right]$$

$$= \sum_{i=1}^N \frac{(y_i - \mu)}{\sigma^2} - \frac{(\mu - \eta)}{\lambda^2} = 0$$

3

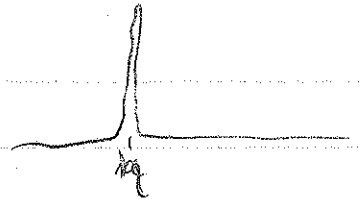
$$\hat{\mu}_{\text{MAP}} = \frac{\left(\sum_{i=1}^N \frac{y_i}{\sigma^2} + \frac{n}{\lambda^2} \right)}{\frac{N}{\sigma^2} + \frac{1}{\lambda^2}}$$

→ Consider, $\lambda \rightarrow \infty$ $p(\mu | n, \lambda^2) = \frac{1}{\sqrt{2\pi}\lambda} e^{-\frac{(\mu-n)^2}{\lambda^2}}$

$p(\mu | n, \lambda)$ is nearly unimodal 

$$\hat{\mu}_{\text{MAP}} \rightarrow \frac{\sum_{i=1}^N \frac{y_i}{\sigma^2} + (\rightarrow 0)}{\frac{N}{\sigma^2} + (\rightarrow 0)} \rightarrow \hat{\mu}_{\text{MLE}} \text{ (nice)}$$

→ Consider $\lambda \rightarrow 0$ $p(\mu | n, \lambda^2)$

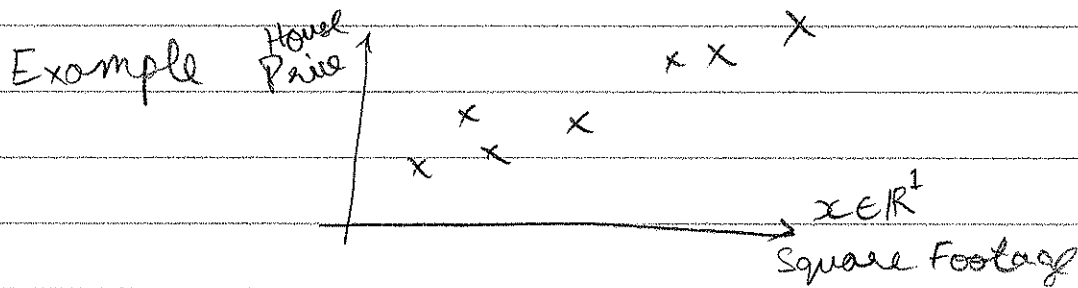


$$\hat{\mu}_{\text{MAP}} \rightarrow \frac{\lambda^2 \sum \frac{y_i}{N} + n}{\lambda^2 \cdot \frac{N}{\sigma^2} + 1}$$

→ $\hat{\mu}_{\text{MAP}} \rightarrow n$ (prior mode)

⑤ Regression

Task: Prediction continuous values $X \rightarrow Y \in \mathbb{R}$



Dataset: $D = \{(\vec{x}_1, y_1), \dots, (\vec{x}_n, y_n)\}$

$\vec{x} \in \mathbb{R}^d$ (Input / features)

$y \in \mathbb{R}$ (Output)

Model:

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_d x_d$$
$$= \begin{bmatrix} w_0 & \dots & w_d \end{bmatrix} \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$= \mathbf{w}^T \vec{x}$$

$$\text{or } \vec{x}^T \mathbf{w}$$

⑥ Least-Square Fitting

Error / Residuals

$$e_i = y_i - \hat{y}_i$$

Ground Truth

prediction

Natural Loss Function $L(y_i, \hat{y}_i) = (y_i - \hat{y}_i)^2$
Squared Loss

(4)

So our objective function / Loss function on the dataset

$$\begin{aligned} \min_{\vec{w} \in \mathbb{R}^{d+1}} L(\vec{w}) &= \frac{1}{n} \sum_{i=1}^n e_i^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \\ &= \frac{1}{n} \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2 \end{aligned}$$

Now we could view $L(\vec{w}) \equiv L(w_0, w_1, \dots, w_d)$
to start chugging along standard calculus

$$\left. \begin{aligned} \frac{\partial L}{\partial w_0} &= 0 \\ \frac{\partial L}{\partial w_1} &= 0 \\ \vdots \\ \frac{\partial L}{\partial w_d} &= 0 \end{aligned} \right\} \text{Solve system of equations}$$

But let's make the notation more compact with matrices

$X_{n \times (d+1)}$
data matrix
"design matrix"

$=$

each training instance is a row
 $n \times (d+1)$

$$Y = \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} \quad \text{all GtS / labels}$$

$$\hat{Y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} \quad \text{all prediction}$$

$$= \begin{bmatrix} \leftarrow \vec{x}_i \rightarrow \end{bmatrix} \begin{bmatrix} w_0 \\ \vdots \\ w_d \end{bmatrix} = Xw$$

$$\text{So } L(\vec{w}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \frac{1}{n} \|Y - \hat{Y}\|_2^2$$

$$= \frac{1}{n} (Y - \hat{Y})^T (Y - \hat{Y})$$

$$= \frac{1}{n} (Y - Xw)^T (Y - Xw)$$

$$= \frac{1}{n} \left[Y^T Y - 2Y^T Xw + w^T X^T Xw \right]$$