# Lecture 11
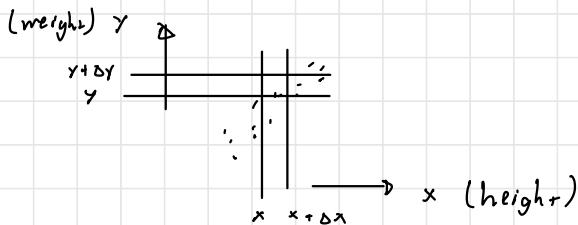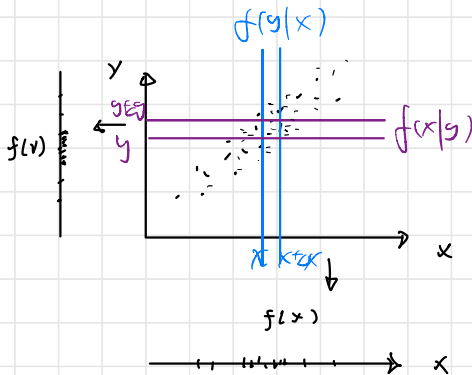
- Probability :

    - $(x, y) \sim f(x, y)$

    - imagine a population of $N (\alpha \infty)$
        equally likely possibilities

    - Population version scatterplot

(weight) Y

Y+∆Y

Y

x   x+∆x

x (height)

○ Counting :
    - joint : $f(x, y)$
    - marginal : $f(x)$, $f(y)$
    - conditional : $f(y|x)$, $f(x|y)$

$f(y|x)$

Y

yEy
y

$f(y)$

$f(x|y)$

x   x+∆x

x

$f(x)$

x

- 3 operations & 1 meta rule

(1) Marginalization:

$$f(x) = \int f(x,y) \, dy$$

$$f(y) = \int f(x,y) \, dx$$
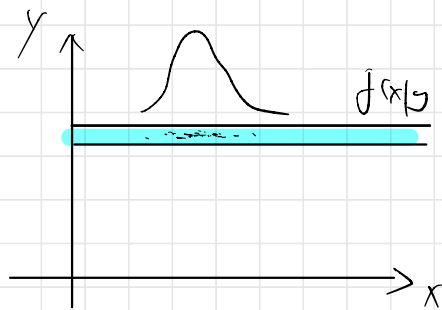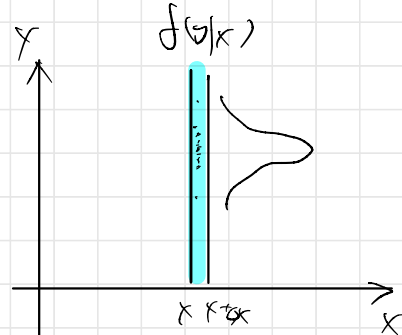
$f(y)$ ←



$f(x)$

(2) Conditioning

$$f(y|x) = \frac{f(x,y)}{f(x)} = \frac{f(x,y)}{\underbrace{\int f(x,y) \, dy}_{\text{normilization}}}$$

$f(y|x)$



$$f(x|y) = \frac{f(x,y)}{f(y)} = \frac{f(x,y)}{\int f(x,y) \, dx}$$

$f(x|y)$



(3) Factorization:

$$f(x,y) = f(x) \, f(y|x) = f(y) \, f(x|y)$$

- Meta Rule: Insert the same condition

Count the same subpopularion

- Example

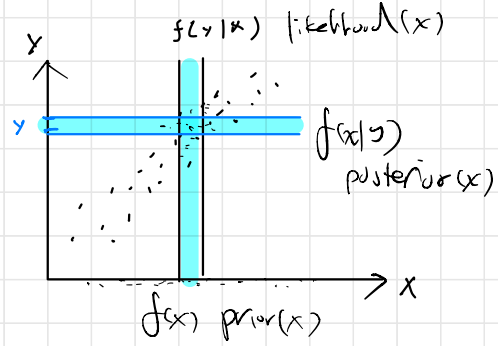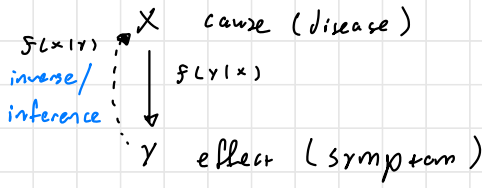(1) $f(x \mid z) = \int f(x, y \mid z) \, dy$     Marginalization

(2) $f(y \mid x, z) = \dfrac{f(x, y \mid z)}{f(x \mid z)}$     Conditioning

(3) $f(x, y \mid z) = f(x \mid z) \, f(y \mid x, z)$     factorization

- Discrete: $\int \longmapsto \Sigma$

- Baye's Rule :

$f(x|y)$
inverse/
inference

$X$ cause (disease)

$\downarrow$ $f(y|x)$

$Y$ effect (symptom)

$$X \sim f(x) \quad \boxed{prior} (x)$$

$$[Y \mid x] \sim f(y|x) \quad \boxed{likelihood} (x)$$

$$[x \mid y] \sim f(x|y) \quad \boxed{posterior} (x)$$

$f(y|x)$ likelihood$(x)$

$Y$

$y$

$f(x|y)$
posterior$(x)$

$X$

$f(x)$ prior$(x)$

Deterministic

$f(y|x)$

point mass at $h(x)$

$Y = h(x)$

$Y$

Solving : $Y = h(x)$
Solution: $x = h^{-1}(y)$

- Baye's Rule:

$$f(x|y) = \frac{f(x,y)}{f(y)} \qquad \text{conditioning}$$

where $x$ is fixed

$$= \frac{f(x,y)}{\int f(x,y)\,dx} \qquad \text{marginalization}$$

$$= \frac{f(x)\,f(y|x)}{\int f(x)\,f(y|x)\,dx} \qquad \text{factorization}$$

- Note: $\quad f(x|y) \propto f(x,y) = f(x)\,f(y|x)$

$\qquad\qquad\qquad\uparrow\qquad\qquad\uparrow\qquad\qquad\downarrow\qquad\quad\downarrow$

$\qquad\qquad\quad$ fixed $\qquad$ fixed $\qquad$ prior $\quad$ likelihood

$\qquad\qquad$ as a function of $x$

- Multivariate Gaussian:

$$X \sim \mathcal{N}(\mu, \Sigma)$$

$$f(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right)$$



$$N \approx \infty$$

$$\Sigma = Q \wedge Q^T$$

$$\begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N}\left(0, \ \Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}\right)$$

Special case when $\Sigma_{21} = \Sigma_{12} = 0$.
Then,

$$f(x_1, x_2) = f(x_1) f(x_2), \quad \text{i.e.} \quad x_1 \perp x_2$$
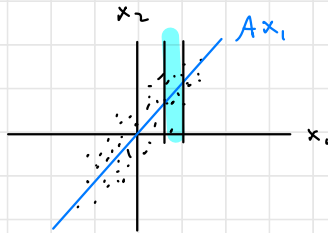
• In general, diagonize

$$\begin{pmatrix} x_1 \\ 2 \end{pmatrix} = \begin{pmatrix} x_1 \\ x_2 - Ax_1 \end{pmatrix} = \begin{pmatrix} I & 0 \\ -A & I \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \sim \mathcal{N}\left( 0, \begin{pmatrix} \Sigma_{11} & 0 \\ 0 & \Sigma_{22} - \Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12} \end{pmatrix} \right)$$

$$A = \Sigma_{21} \Sigma_{11}^{-1}$$

$$X_2 = Ax_1 + \varepsilon, \qquad \varepsilon \perp x_1$$

$$[x_2 \mid x_1] \sim \mathcal{N}( Ax_1, \ \Sigma_{22} - \dots )$$

<span style="color:blue">↑<br>fixed</span>



○ Change Notation:

$$\begin{pmatrix} x \\ y \end{pmatrix} \sim \mathcal{N}\left( 0, \begin{pmatrix} \Sigma_{xx} & \Sigma_{xy} \\ \Sigma_{yx} & \Sigma_{yy} \end{pmatrix} \right)$$

- population - version regression

| | $x_i^T$ | $y_i^T$ |
|---|---|---|
| 1 ⋮ i ⋮ N | X | Y |

$$\hat{\beta}_{LS} = ( X^T X )^{-1} (X^T y )$$

$$\frac{1}{N} X^T X =$$



$$= \frac{1}{N} \sum_{i=1}^{N} x_i x_i^T \overset{1\text{dim}}{=} \frac{1}{N} \sum x_i^2 = \frac{1}{N} |\vec{x}|^2$$

$$= \mathbb{E}\left((x-\mu)(x-\mu)^T\right) = Var(x) = \Sigma_{xx}$$

**Similarly**

$$\frac{1}{N} X^T Y = \Sigma_{xy} \overset{1\text{dim}}{=} \frac{1}{N} \langle \vec{x}, \vec{y} \rangle$$

**So** $\quad \hat{\beta} = \Sigma_{xx}^{-1} \Sigma_{xy} = A^T$

1 dim



Scatterplot



$y = \hat{\beta} x$

$(x_i, y_i)$

2 dim

Vector plot



$\Sigma_{YY}$

$\Sigma_{YY} - \Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$

$\Sigma_{YX} \Sigma_{XX}^{-1} \Sigma_{XY}$

$\Sigma_{XX}$

**Correlation**

$$\rho = Corr(x,y) = \frac{Cov(x,y)}{\sqrt{Var(x)}\sqrt{Var(y)}} = Cov\left(\frac{x-\mu_x}{\sigma_x}, \frac{y-\mu_y}{\sigma_y}\right)$$

$$= \frac{\frac{1}{N}\langle \vec{x}, \vec{y} \rangle}{\sqrt{\frac{1}{N}|\vec{x}|^2}\sqrt{\frac{1}{N}|\vec{y}|^2}} = \cos\theta$$

$$\frac{x-\mu_x}{\sigma_x} \rightarrow \begin{pmatrix} x \\ y \end{pmatrix} \sim W\left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} 1 & \rho \\ \rho & 1 \end{pmatrix}\right)$$

$$\frac{y-\mu_y}{\sigma_y} \rightarrow$$

**Regression**

$$[y|x] \sim N(\rho x, (1-\rho^2))$$

$$1-\rho^2 = \sin^2\theta = \frac{|\vec{e}|^2}{|\vec{y}|^2} = \frac{Var(e)}{Var(y)}$$

Son height $\quad$ father height $\quad < 1$

Regression towards the mean

original meaning of "regression"

- **Bayesian Regression:**

Data

| | $p$ |
|---|---|
| $X$ | $Y$ |

$n$

$$[Y \mid X, \beta] \sim N(X\beta, \sigma^2 I_n)$$

$$\overset{P(\beta)}{\text{prior:}} \quad \beta \sim N(0, \tau^2 I_p)$$

(as if $\beta$ is sampled from population of $N$ possibilities)

posterior:
$$P(\beta \mid X, Y) \propto \underset{\text{prior}(\beta)}{P(\beta)} \ \underset{\text{likelihood}(\beta)}{P(Y \mid X, \beta)}$$

Note: Gauss assumed $\tau^2 = \infty$, so we know nothing about $\beta$ (non-informative prior).

- **What about finite $\tau^2$**

$$P(\beta \mid X, Y) \propto \frac{1}{(2\pi)^{\frac{p}{2}} (\tau^2)^{p/2}} \exp\left(-\frac{1}{2\tau^2} \beta^T \beta\right)$$
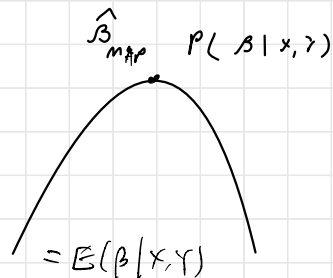
$$\begin{array}{c}[Y \mid X, \beta] \\ \text{density}\end{array} \longrightarrow \frac{1}{(2\pi)^{n/2} (\sigma^2)^{\frac{n}{2}}} \exp\left(-\frac{1}{2\sigma^2}(Y - X\beta)^T (Y - X\beta)\right)$$

$$\propto \exp\left(-\frac{1}{2}\left(\frac{1}{\sigma^2}|Y - X\beta|^2 + \frac{1}{\tau^2}|\beta|^2\right)\right)$$

Ridge - Regression.

- **Maximum A Posteriori: (MAP)**

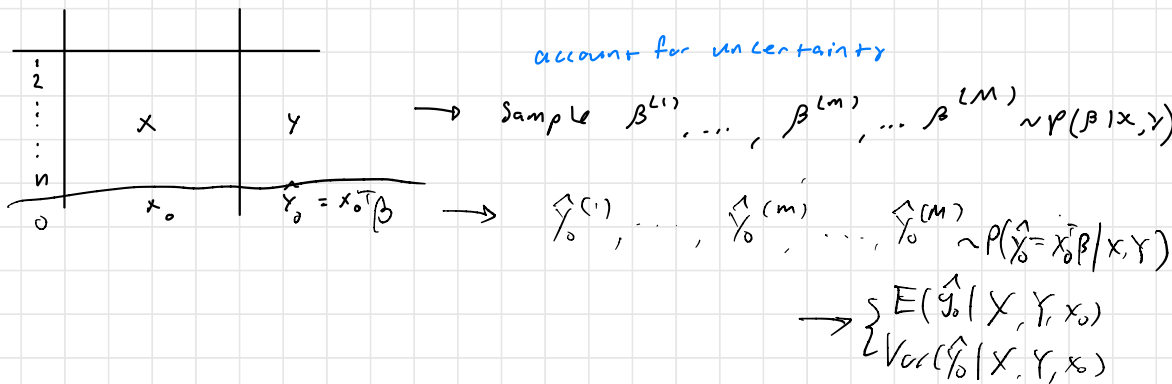$$\min_{\beta} \frac{1}{\sigma^2}|Y - X\beta|^2 + \frac{1}{\tau^2}|\beta|^2$$

$$= \min_{\beta} |Y - X\beta|^2 + \underset{\uparrow \ \lambda > 0}{\frac{\sigma^2}{\tau^2}}|\beta|^2$$

$$\hat{\beta}_{MAP} \quad P(\beta \mid X, Y)$$

$$= \mathbb{E}(\beta \mid X, Y)$$

- $[\beta \mid X, Y] \sim p(\beta \mid X, Y)$

$\downarrow$

basis for inference / prediction

$\downarrow$

Does not overfit

account for uncertainty



$\rightarrow$ Sample $\beta^{(1)}, \ldots, \beta^{(m)}, \ldots \beta^{(M)} \sim p(\beta \mid X, Y)$

$x_0 \qquad \hat{Y}_0 = x_0^T \beta \rightarrow \hat{Y}_0^{(1)}, \ldots, \hat{Y}_0^{(m)}, \ldots, \hat{Y}_0^{(M)} \sim p(\hat{Y}_0 = x_0^T \beta \mid X, Y)$

$\rightarrow \begin{cases} E(\hat{Y}_0 \mid X, Y, x_0) \\ Var(\hat{Y}_0 \mid X, Y, x_0) \end{cases}$

Can be derived directly

population of $\beta \rightarrow$ population of $\begin{pmatrix} Y \\ \hat{Y}_0 \end{pmatrix} \sim N\left( 0, \qquad \right)$
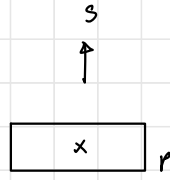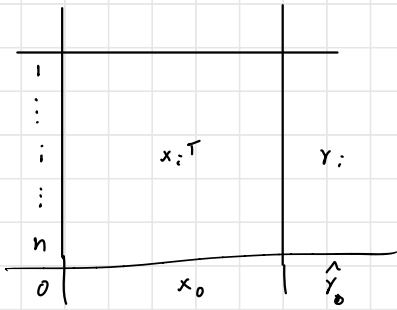
- Bayesian / Frequentist controversy

  - example: prior, speed of light $c \sim$ prior $(c)$

  - $c \sim p(c)$
  - feasible of repeated sampling ?
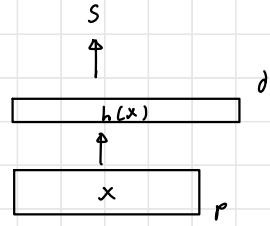
    $c_1, c_2, \ldots, c_m \overset{iid}{\sim} p(c)$

    $\downarrow$

    frequency $\rightarrow$ prob

For the diagram elements:

Table:
| | $x_i^T$ | $Y_i$ |
| 1 | | |
| : | | |
| i | | |
| : | | |
| n | | |
| 0 | $x_0$ | $\hat{Y}_0$ |

$$s = f(x) = x^T \beta$$
$$\beta \sim N(0, \sigma^2 I_\rho)$$
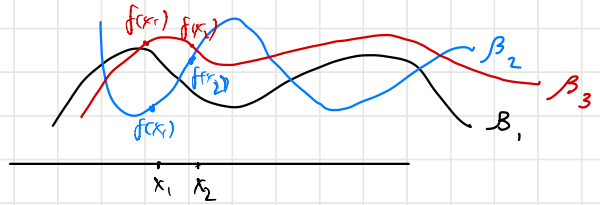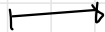$$(Ridge)$$

$$s = f(x) = h(x)^T \beta$$
$$\beta \sim N(0, \sigma^2 I_d)$$
$$(kernel)$$

- For simplicity assume $X$ is $1d$

Population of $\beta \longrightarrow$ population of $f(x) = h(x)^T \beta$

$\beta$ pop $\longmapsto$

$f(x_1) \quad f(x_2)$
$f(x_1)$
$f(x_2)$
$\beta_2$
$\beta_3$
$\beta_1$
$x_1 \quad x_2$

So me say $f(x) = h(x)^T \beta \sim$ Gaussian Process

So $\text{Cov}(f(x_1), f(x_2)) = \text{Cov}(h(x_1)^T B, \ h(x_2)^T B)$

Remember $\mu = 0$

$$= \mathbb{E}(h(x_1)^T B B^T h(x_2)) \quad , B \text{ is random}$$
$$x_1 \& x_2 \text{ is fixed}$$

$$= h(x_1)^T \mathbb{E}(B B^T) h(x_2)$$

$$= h(x_1)^T \sigma^2 I_d \ h(x_2)$$

$$= \sigma^2 \langle h(x_1), h(x_2) \rangle$$

$$= \sigma^2 k(x_1, x_2)$$

So $f(x) \sim GP(0, \sigma^2 k)$

| 1 | | $f(x_1)$ | $y_1$ |
|---|---|---|---|
| $\vdots$ | | | |
| $i$ | $x_i^T$ | $f(x_i)$ | $y_c = f(x_i) + \varepsilon_i$ |
| $\vdots$ | | | |
| $n$ | | $f(x_n)$ | $y_n$ |
| $0$ | $x_0$ | $f(x_0)$ | |

$$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_i \\ \vdots \\ \varepsilon_n \end{pmatrix} \sim N(0, \sigma^2 I)$$

$$\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_i) \\ f(x_n) \\ \hline f(x_0) \end{pmatrix} \sim N \left( 0, \ \begin{pmatrix} & \sigma^2 k_{ij} & \sigma k_{i0} \\ & & \\ \hline & \sigma^2 k_{0j} & \sigma^2 k_{00} \end{pmatrix} \right)$$

$$k_{ij} = k(x_i, x_j)$$

$$\begin{pmatrix} y_1 \\ \vdots \\ y_i \\ \vdots \\ y_n \\ \hat{y}_0 = f(x_0) \end{pmatrix} \sim N \left( 0, \quad \begin{array}{|c|c|} \hline & \tau^2 k_0^T \\ \tau^2 k + \delta^2 I_n & \\ \hline \tau^2 k_0 & \sigma^2 k_{00} \\ \hline \end{array} \right)$$

$$S_0 \left[ \hat{y}_0 = f(x_0) \mid Y, x \right] \sim N \left( \tau^2 k_0 \left( \tau^2 k + \delta^2 I_n \right)^{-1} Y, \right.$$

$$\left. \tau^2 k_{00} - \tau^2 k_0 \left( \tau^2 k + \sigma^2 I_n \right)^{-1} \tau^2 k_0^T \right)$$

Kernel Regression

$$k_0 \left( k + \frac{\sigma^2}{\tau^2} I_n \right)^{-1} Y$$

$$= k_0 \left( K + \lambda I_n \right)^{-1} Y$$

$$= k_0 \, C$$

$$= \boxed{\quad k_0 \quad} \; \boxed{\begin{array}{c} c \end{array}}$$

$$f(x_0) = \sum_{i=1}^{n} c_i \, k(x_i, x_0)$$