## Bayesian Hypothesis Testing

Consider $\qquad H_0: \theta = \theta_0$

$\qquad$ vs. $\quad H_1: \theta \neq \theta_0$

Technique : $\quad$ Put a prior on $\theta$ <u>and</u> $H_0$, then

$\qquad$ compute : $\quad \mathbb{P}(H_0 \mid \vec{X})$

$\qquad\qquad\qquad\qquad$ $\curvearrowright$ this is our observed

$\qquad\qquad\qquad\qquad\qquad$ data.

Ex: $\quad$ Put the prior $\quad P(H_0) = \frac{1}{2}$, $\quad P(H_1) = \frac{1}{2}$.

Then $\quad$ compute $\qquad \mathbb{P}(H_0 \mid \vec{X})$

$$\mathbb{P}(H_0 \mid \vec{X}) = \frac{f(\vec{X}, H_0)}{f(\vec{X})}$$

$$= \frac{f(\vec{X} \mid H_0) \cdot P(H_0)}{f(\vec{X} \mid H_0) P(H_0) + f(\vec{X} \mid H_1) P(H_1)}$$

in our case,

$$P(H_0) = P(H_1)$$

$$= \frac{f(\vec{X} \mid \theta_0)}{f(\vec{X} \mid \theta_0) + f(\vec{X} \mid H_1)}$$

$$= \frac{f(\vec{X} \mid \theta_0)}{f(\vec{X} \mid \theta_0) + \int f(\vec{X} \mid \theta) f(\theta) d\theta}$$

$$= \frac{\mathcal{L}(\theta_0)}{\mathcal{L}(\theta_0) + \int \mathcal{L}(\theta) f(\theta) d\theta}$$

<u>Notes</u> — prior $f$ can have a large influence on
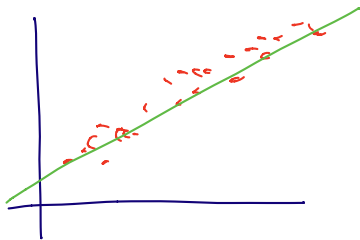
$H_0 | \vec{x}$

— Improper priors are <u>not</u> allowed

— $\mathbb{P}(H_0 | \vec{x})$ is the probability that $H_0$ is true given $\vec{x}$ → this does not tell us when to reject the null hypothesis. When do we reject? When do we retain? We need more detailed analysis.
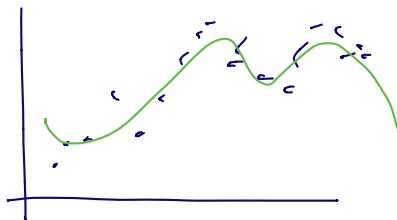
Read 11.9 in All of Stats for more strengths/weaknesses.
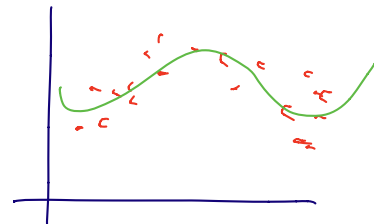
<u>Regression</u>  (standard linear regression)

<u>Goal</u>: Fit noisy data using a curve.



$y = ax + b$

parametric

$y = a_0 + a_1 x + a_2 x^2 + \ldots + a_n x^n$

parametric

$y | \vec{x}, \vec{y} \sim GP(m, k)$

non-parametric

Denote:  $Y \sim$ random variable , response variable

$X \sim$ covariate, predictor, feature

Regression:  $r(x) = \mathbb{E}(Y | \vec{X} = \vec{x}) = \int y \, f(y | \vec{x}) \, dy$

Goal:  Given data $(X_1, Y_1) \ldots (X_n, Y_n) \sim F_{XY}$ , estimate $r(x)$.

model ↗

# Basic Linear Regression

Model : $\quad r(x) = \beta_0 + \beta_1 x \qquad$ "simple linear regression model"

Observe some data : $\quad X_i, Y_i$

Assumption : $\quad Var(Y \mid X = x) = \sigma^2$

$\quad\quad\quad \hookrightarrow Y_i = \beta_0 + \beta_1 X_i + \epsilon_i$

$$\mathbb{E}(\epsilon_i \mid X_i) = 0$$

$$Var(\epsilon_i \mid X_i) = \sigma^2.$$

Given this data for the statistical model, find estimators for the unknown coefficients $\beta_0, \beta_1 \to \hat{\beta}_0, \hat{\beta}_1$.
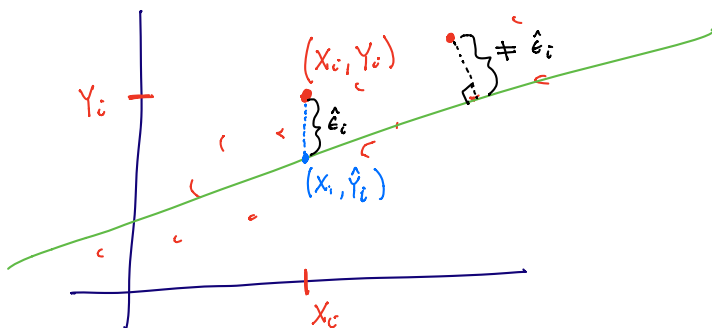
Fitted line : $\quad \hat{r}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$

Predicted values : $\quad \hat{Y}_i = \hat{r}(x_i)$

Residuals : $\quad \hat{\epsilon}_i = Y_i - \hat{Y}_i = Y_i - \left(\hat{\beta}_0 + \hat{\beta}_1 X_i\right)$

Residual sum of squares : $\quad RSS = \sum \hat{\epsilon}_i^2 \quad \leftarrow$ <span style="color:red">only one such metric to determine how well $\hat{r}$ fits the data.</span>



## Definition :

Least squares estimates : $\hat{\beta}_0, \hat{\beta}_1$

minimize $RSS = \sum_{i=1}^{n} \hat{\epsilon}_i^2$.

$\hat{\beta}_0, \hat{\beta}_1$ can be found using calculus, linear algebra, statistics, etc.

③

Thm: $\hat{\beta}_1 = \dfrac{\sum_i (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_i (X_i - \bar{X})^2}$ $\Bigg\}$ $\approx \dfrac{cov(X,Y)}{var(X)}$

and $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$

$$= \dfrac{\sigma_X \sigma_Y \rho_{XY}}{\sigma_X^2}$$

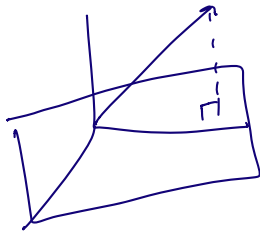$$= \dfrac{\sigma_Y}{\sigma_X} \rho_{XY}$$

Unbiased estimate of $\sigma^2$ :

$$\hat{\sigma}^2 = \dfrac{1}{n-2} \sum \hat{\epsilon}_i^2 \;.$$

To find $\hat{\beta}_0, \hat{\beta}_1$ using calculus :    solve $\dfrac{\partial}{\partial \beta_0} RSS = 0$

$$\dfrac{\partial}{\partial \beta_1} RSS = 0$$

To solve using linear algebra :



compute the orthogonal projection
of $Y_i$ onto span $\{\vec{1}, \vec{X}\}$

## Least Squares and Maximum Likelihood Estimators

Add assumption that $\epsilon_i | X_i \sim N(0, \sigma^2)$

$$\Rightarrow Y_i | X_i \sim N(\mu_i, \sigma^2)$$

$$\hookrightarrow \mu_i = \beta_0 + \beta_1 X_i$$

Write down the likelihood function :

$$\mathcal{L} \;\propto\; \prod f(X_i, Y_i) \;=\; \underbrace{\prod f(X_i)}_{\mathcal{L}_1} \cdot \underbrace{f(Y_i | X_i)}_{}$$

$$= \mathcal{L}_1 \quad \times \quad \mathcal{L}_2(\beta_0, \beta_1, \sigma^2) \quad \boxed{4}$$

$\mathcal{L}_1$ does not depend on any parameters.

$\mathcal{L}_2$ is known as the __conditional likelihood__ and contains all the parameters.

$$\mathcal{L}_2(\beta_0, \beta_1, \sigma^2) \propto \frac{1}{\sigma^n} e^{-\frac{1}{2\sigma^2} \sum_i (Y_i - \beta_0 - \beta_1 X_i)^2}$$

$$\Rightarrow \log \mathcal{L}_2 = \ell_2 = -n \log \sigma - \frac{1}{2\sigma^2} \underbrace{\sum (Y_i - \beta_0 - \beta_1 X_i)^2}$$

__Thm__ If $\epsilon_i | X_i \sim N(0, \sigma^2)$, then the MLE
for $\beta_0, \beta_1$ is __the same__ as the least squares estimate,
and $\hat{\sigma}^2_{MLE} = \frac{1}{n} \underbrace{\sum \hat{\epsilon}_i^2}$
biased estimator.

__Properties of these least squares estimators__

__Thm__: Define $\hat{\beta} = \begin{pmatrix} \hat{\beta}_0 \\ \hat{\beta}_1 \end{pmatrix}$

then $\mathbb{E}\left( \hat{\beta} | \vec{X} \right) = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$

$$\text{Var}\left( \hat{\beta} | \vec{X} \right) = \frac{\sigma^2}{n S_x^2} \begin{pmatrix} \frac{1}{n} \sum X_i^2 & -\bar{X} \\ -\bar{X} & 1 \end{pmatrix}$$

where $S_x^2 = \frac{1}{n} \sum (X_i - \bar{X})^2$

The estimated standard errors are ( sqrt of diagonals of covariance matrix )

$$\hat{se}(\hat{\beta}_0) = \frac{\hat{\sigma}}{\sqrt{n}\, s_x} \sqrt{\frac{1}{n} \sum x_i^2}$$

$$\hat{se}(\hat{\beta}_1) = \frac{\hat{\sigma}}{\sqrt{n}\, s_x}$$

**Thm** These estimators are

- consistent

- asymptotically normal

- and therefore we can apply the Wald test,

  e.g. $H_0: \beta_1 = 0$    vs.    $H_1: \beta_1 \neq 0$.

## Prediction

Setup : Have the data $X_1, Y_1 \dots X_n, Y_n$ and

estimate $\hat{r}(x) = \hat{\beta}_0 + \hat{\beta}_1 x$   using   least squares.

Observe ( or pick ) a new covariate $X = x_*$, and

we want to predict $Y_*$.

**Estimate** $\hat{Y}_* = \hat{\beta}_0 + \hat{\beta}_1 x_*$,

$$Var(\hat{Y}_*) = Var(\hat{\beta}_0 + \hat{\beta}_1 x_*)$$

$$= Var(\hat{\beta}_0) + x_*^2 \, Var(\hat{\beta}_1) + 2x_* \, Cov(\hat{\beta}_0, \hat{\beta}_1)$$

and $\hat{se}(\hat{Y}_*) = \sqrt{Var(\hat{Y}_*)}$ using $\hat{\sigma}^2$.

How about a $(1-\alpha)$ confidence interval for $Y_*$?

Mainly, $\hat{Y}_* \pm z_{\alpha/2}\, \hat{se}(\hat{Y}_*)$ is a $1-\alpha$ confidence interval,

<span style="color:red">but this is <u>incorrect</u>.</span> ( See Thm 13.11, do exercise 10)

The idea behind the mistake:

the above confidence interval is only correct if we never observed the independent noise $\epsilon_i$, i.e., in the real world we observe $Y_* = \beta_0 + \beta_1 X_* + \epsilon$.

# Multiple Regression

Setup: $\vec{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ , covariats $X = \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1k} \\ X_{21} & & & \vdots \\ \vdots & & & \vdots \\ X_{n1} & & & X_{nk} \end{pmatrix}$

$$2^{nd} \text{ covariate.}$$

$$\vec{\beta} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad , \quad \vec{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

Model: $\vec{Y} = X\vec{\beta} + \vec{\epsilon}$

See Thm 13.13 for the least squares solution, same as from linear algebra class.
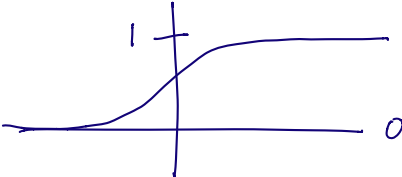
# Logistic Regression

Change the model: Imagine that $Y_i \in \{0,1\}$,

i.e. $P(Y_i | X_i) = p_i$

We want to model $p_i$, not $Y_i$.

Choose a particular parametric form:

$$p_i = p_i(\beta_0, .., \beta_k) = \mathbb{P}(Y_i = 1 \mid X_i)$$

$$= \frac{e^{\beta_0 + \sum_i^k \beta_j x_{ij}}}{1 + e^{\beta_0 + \sum_i^k \beta_j x_{ij}}}$$

The logistic function: $\dfrac{e^x}{1 + e^x}$



$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right)$$

Since $Y_i \mid X_i \sim \text{Bernoulli}(p_i)$

The conditional likelihood is:

$$\mathcal{L}(\vec{\beta}) = \prod_i p_i(\vec{\beta})^{Y_i} \left(1 - p_i(\vec{\beta})\right)^{1-Y_i}$$

$\mathcal{L}$ must be maximized numerically.

## Multivariate Models

Random vector $\vec{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_k \end{pmatrix}$, mean $\mathbb{E}(\vec{X}) = \begin{pmatrix} \mathbb{E}(X_1) \\ \mathbb{E}(X_2) \\ \vdots \\ \mathbb{E}(X_k) \end{pmatrix} = \vec{\mu}$

8

Covariance matrix

$$C = \begin{pmatrix} cov(X_1, X_1) & cov(X_1, X_2) & \cdots & cov(X_1, X_k) \\ cov(X_2, X_1) & & & \vdots \\ cov(X_3, X_3) & & & \vdots \\ \vdots & & & \\ cov(X_k, X_1) & & & cov(X_k, X_k) \end{pmatrix}$$

$\underbrace{\qquad\qquad\qquad\qquad}_{}$

$k \times k$ matrix.

$C^{-1}$ is known as the <u>precision matrix</u>,

Furthermore :  — $C$ is symmetric positive definite

— eigenvalues of $C$ are all positive.

<u>Thm</u>    Let    $\vec{a} \in \mathbb{R}^k$  be  a  constant vector,

then    ①   $\mathbb{E}(\vec{a}^T \vec{X}) = \vec{a}^T \vec{\mu}$

②  $Var(\vec{a}^T \vec{X}) = \vec{a}^T C \vec{a}$

Let  $A \in \mathbb{R}^{n \times k}$  constant matrix  ,  then

①  $\mathbb{E}(A\vec{X}) = A\vec{\mu}$

②  $Var(A\vec{X}) = A C A^T$  } another covariance matrix

Next  consider  we have  samples

$$X_{11}, X_{12}, \ldots, X_{1n}$$
$$\vdots$$
$$X_{k1}, \ldots \ldots - X_{kn}$$

The sample mean is then

$$\bar{X} = \begin{pmatrix} \bar{X}_1 = \frac{1}{n}\sum_j X_{1j} \\ \vdots \\ \bar{X}_k = \frac{1}{n}\sum_j X_{kn} \end{pmatrix}$$

The sample variance matrix:

$$S = \begin{pmatrix} S_{11} & S_{12} & \cdots \\ & \ddots & \\ & & S_{kk} \end{pmatrix}$$

where
$$S_{ij} = \frac{1}{n-1}\sum_{\ell=1}^{n}(X_{i\ell} - \bar{X}_i)(X_{j\ell} - \bar{X}_j) \left.\right\} \begin{array}{l} \text{unbiased} \\ \text{estimate of} \\ \text{Cov}(X_i, X_j). \end{array}$$

$$\mathbb{E}(\bar{X}) = \vec{\mu} \qquad , \qquad \mathbb{E}(S) = C .$$

And since the correlation of two variables is

$$\rho(X_i, X_j) = \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}(X_i)\,\text{Var}(X_j)}} \qquad , \qquad \text{the plugin}$$

estimator for the correlation is:

$$\hat{\rho}_{ij} = \frac{S_{ij}}{S_{ii}\,S_{jj}}$$

Example   Multivariate Normal

$\vec{X} \in \mathbb{R}^k \sim N(\vec{\mu}, C)$   if   its density is

$$f(x_1, \ldots, x_k; \vec{\mu}, C) = \frac{1}{(2\pi)^{k/2}} \frac{1}{\sqrt{\det C}} e^{-\frac{1}{2}(\vec{X} - \vec{\mu})^T C^{-1}(\vec{X} - \vec{\mu})}$$

$\vec{\mu} \in \mathbb{R}^k$

$C \in \mathbb{R}^{k \times k}$

$\Rightarrow \mathbb{E}(\vec{X}) = \vec{\mu}$

$\text{Var}(\vec{X}) = C$.

Thm   Let $\vec{Z} \sim N(\vec{0}, I)$, and $C$ be a spd matrix.

① Let $C^{1/2}$ be such that $C^{1/2} \cdot C^{1/2} = C$, then

$$\vec{X} = \vec{\mu} + C^{1/2} \vec{Z} \sim N(\vec{\mu}, C).$$

② $C^{-1/2}(\vec{X} - \vec{\mu}) \sim N(\vec{0}, I)$.

③ $\vec{a}^T \vec{X} \sim N(\vec{a}^T \vec{\mu}, \vec{a}^T C \vec{a})$

④ $V = (\vec{X} - \vec{\mu})^T C^{-1} (\vec{X} - \vec{\mu})$, then $V \sim \chi_k^2$

$= \vec{Z}^T \vec{Z}$.

Thm For multivariate dat $\begin{pmatrix} X_{11} \\ X_{21} \\ \vdots \\ X_{k1} \end{pmatrix} \cdots \begin{pmatrix} X_{1n} \\ \vdots \\ X_{kn} \end{pmatrix}$, the log-likelihood

$$\ell(\vec{\mu}, C) \propto \frac{-n}{2}(\vec{\bar{x}} - \vec{\mu})^T C^{-1}(\vec{\bar{x}} - \vec{\mu}) - \frac{n}{2} \text{tr}(C^{-1} S) - \frac{n}{2} \log \det C$$

[11]

Recall that $\text{tr}(A) = \sum A_{ii}$,

$$S = \text{as before, the sample covariance matrix.}$$

The MLE's are $\hat{\mu} = \overline{x}$, $\hat{C} = \frac{n-1}{n} S$

# Gaussian Processes

Simplest interpretation: the extension of random vectors to random functions

Best single source reference:

Rasmussen & Williams, Gaussian Processes for Machine Learning

We say that $f \sim GP(m, k)$ is a Gaussian process with mean $m$ and covariance function $k$.

$$\Rightarrow \quad \mathbb{E}(f(x)) = m(x)$$
$$\text{Cov}(f(x), f(x')) = k(x, x')$$

$$\Rightarrow \quad \mathbb{E}\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} = \begin{pmatrix} m(x_1) \\ m(x_2) \\ \vdots \\ m(x_n) \end{pmatrix}$$

$$\text{Cov}\begin{pmatrix} f(x_1) \\ \vdots \\ f(x_n) \end{pmatrix} = K \quad \text{with entries} \quad k(x_i, x_j)$$

$$\Leftrightarrow \quad k(x, x') = \mathbb{E}\left( (f(x) - m(x))(f(x') - m(x')) \right)$$
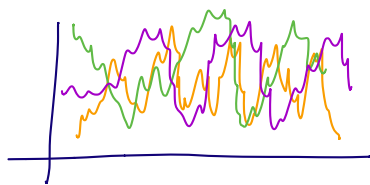
Example    covariance functions:

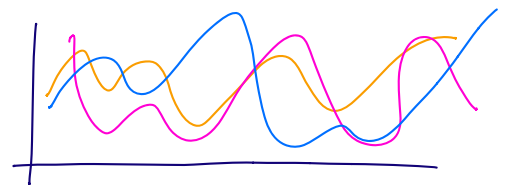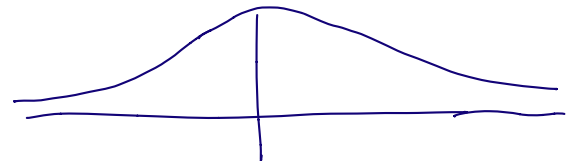$$k(x, x') = A e^{-(x-x')^2/b}$$

$$k(x, x') = B e^{-|x-x'|/c} \left.\right\} \text{ one function from the Matérn family of covariance kernels.}$$

Graphically

$$k(x, x') = e^{-(x-x')^2/.0001}$$



$r = |x - x'|$



$$k(x, x') = e^{-(x-x')^2/10000}$$





$k(x, x') = \delta(x-x') = 1$ if $x = x'$
$0$ otherwise

$w(t)$



Discontinuous everywhere.

Brownian Motion  $B(t) = \int_0^t w(\tau) \, d\tau$

$\hookleftarrow GP(0, \delta)$