

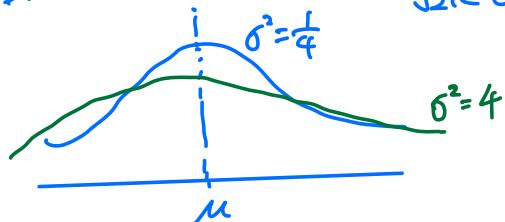
CS136 lecture 6 Gaussian model

Outline

1. Univariate Gaussian distribution
2. MLE for Gaussians
3. Biased & unbiased estimators
4. Special properties of Gaussian
5. Covariance and covariance Matrices

Univariate Gaussian distribution

$$\text{PDF: Normal PDF}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



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

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

$$\Rightarrow \int e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \sqrt{2\pi}\sigma$$

Computing its mean and variance

Mean:

$$\begin{aligned} E(X) &= \int_{-\infty}^{\infty} x \cdot \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} x e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \\ &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} X e^{-\frac{1}{2\sigma^2}(X-\mu)^2} dx \quad * \end{aligned}$$

$$\text{let } t = \frac{x-\mu}{\sigma} \Rightarrow x = \sigma t + \mu$$

$$dt = \frac{1}{\sigma} dx \Rightarrow dx = \sigma dt$$

$$\begin{aligned} * &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} (\sigma t + \mu) \cdot t e^{-\frac{t^2}{2}} dt \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} (\sigma t + \mu) \cdot e^{-\frac{t^2}{2}} dt \\
&= \frac{1}{\sqrt{2\pi}} \left(\int_{-\infty}^{+\infty} \sigma t \cdot e^{-\frac{t^2}{2}} dt + \int_{-\infty}^{+\infty} \mu e^{-\frac{t^2}{2}} dt \right) \\
&= \frac{1}{\sqrt{2\pi}} \left(\sigma \underbrace{\int_{-\infty}^{+\infty} t e^{-\frac{t^2}{2}} dt}_{A} + \mu \underbrace{\int_{-\infty}^{+\infty} e^{-\frac{t^2}{2}} dt}_{B} \right) \\
&= \frac{1}{\sqrt{2\pi}} (\sigma \cdot 0 + \underbrace{\sqrt{2\pi} \mu}_{A}) \\
&= \mu
\end{aligned}$$

<p>A:</p> $f(t) = \int_{-\infty}^{+\infty} t e^{-\frac{t^2}{2}} dt$ $f(-t) = - \int_{-\infty}^{+\infty} t e^{-\frac{t^2}{2}} dt$ $-f(t) = - \int_{-\infty}^{+\infty} t e^{-\frac{t^2}{2}} dt$ <p>odd function: $f(-x) = -f(x)$</p> $\int_{-\infty}^{+\infty} \text{odd function} = 0$ 	<p>B:</p> $\int e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx = \sqrt{2\pi} \sigma$ <p>here. $\mu=0, \sigma=1$</p> $\Rightarrow \int e^{-\frac{t^2}{2}} dt = \sqrt{2\pi}$
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------------------------------------------------------------------------------------

Variance

$$\begin{aligned}
\text{Var}(X) &= E(X^2) - \overline{E(X)}^2 \\
E(X^2) &= \int_{-\infty}^{+\infty} x^2 \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} dx \\
&= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} x^2 e^{-\frac{1}{2} \cdot \frac{(x-\mu)^2}{\sigma^2}} dx
\end{aligned}$$

$$\begin{aligned}
\text{let } t &= \frac{x-\mu}{\sigma}, \quad \Rightarrow x = \sigma t + \mu \\
dt &= \frac{1}{\sigma} dx \quad dx = \sigma dt
\end{aligned}$$

$$\begin{aligned}
\text{**} &\Rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} (\sigma t + \mu)^2 \cdot e^{-\frac{t^2}{2}} dt \\
&= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{+\infty} (\sigma^2 t^2 + 2\mu\sigma t + \mu^2) e^{-\frac{t^2}{2}} dt \\
A: \quad &\int_{-\infty}^{+\infty} \sigma^2 t^2 e^{-\frac{t^2}{2}} dt = \sigma^2 \int_{-\infty}^{+\infty} t^2 e^{-\frac{t^2}{2}} dt = \sqrt{2\pi} \sigma^2
\end{aligned}$$

Δ

$$B: \int_{-\infty}^{+\infty} (2\mu\sigma t e^{-\frac{t^2}{2}} dt) = 2\mu\sigma \left(\int_{-\infty}^{+\infty} e^{-\frac{t^2}{2}} dt \right) = 0$$

$$C: \int_{-\infty}^{+\infty} \mu^2 e^{-\frac{t^2}{2}} dt = \mu^2 \int_{-\infty}^{+\infty} e^{-\frac{t^2}{2}} dt = \sqrt{2\pi} \mu^2$$

$$\Rightarrow \frac{1}{\sqrt{2\pi}} (\sqrt{2\pi}\sigma^2 + \sqrt{2\pi}\mu^2) = \sigma^2 + \mu^2$$

$$\Rightarrow E(X^2) = \sigma^2 + \mu^2$$

Besides, $E(X)^2 = \mu^2$

$$\therefore \text{Var}(X) = \sigma^2$$

MLE for gaussians / biased & unbiased estimators

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

$$\text{MLE: } f(y; \mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

for μ :

$$\begin{aligned} f(y; \mu) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(x_i - \mu)^2}{2\sigma^2} \\ &= \sum_{i=1}^n -(\log \sqrt{2\pi} + \log \sigma) - \frac{(x_i - \mu)^2}{2\sigma^2} \end{aligned}$$

$$\frac{\partial f(y; \mu)}{\partial \mu} = \cancel{\sum_{i=1}^n} \sum_{i=1}^n (x_i - \mu) \stackrel{\text{set}}{=} 0$$

$$\Rightarrow \sum_{i=1}^n x_i - n\mu = 0$$

$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$E(\hat{\mu}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} n \cdot \mu$$

$$= \mu \quad \rightarrow \text{unbiased estimator}$$

$$\begin{aligned}
 f(y; \sigma) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \\
 &= \prod_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{(x_i - \mu)^2}{2\sigma^2} \\
 &= \prod_{i=1}^n -(\log \sqrt{2\pi} + \log \sigma) - \frac{(x_i - \mu)^2}{2\sigma^2} \\
 &= -\sum_{i=1}^n (\log \sigma + \frac{(x_i - \mu)^2}{2}) + \frac{1}{\sigma^2}
 \end{aligned}$$

$$\frac{\partial f(y; \sigma)}{\partial \sigma} = -\sum_{i=1}^n \left(\frac{1}{\sigma} - (x_i - \mu)^2 \cdot \frac{1}{\sigma^3} \right) \stackrel{\text{set}}{=} 0$$

$$\times \sigma \Rightarrow \sum_{i=1}^n (\sigma^2 - (x_i - \mu)^2) = 0$$

$$n\sigma^2 - \sum_{i=1}^n (x_i - \mu)^2 = 0$$

$$\boxed{\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

← biased estimation

$$\begin{aligned}
 E(\sigma^2) &= E\left(\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2\right) \\
 &= \frac{1}{n} E\left(\sum_{i=1}^n (x_i - \hat{\mu})^2\right) \\
 &= \frac{1}{n} E\left[\sum_{i=1}^n x_i^2 - 2 \sum_{i=1}^n x_i \hat{\mu} + \sum_{i=1}^n \hat{\mu}^2\right] \\
 &= E\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - 2 \underbrace{\frac{1}{n} \sum_{i=1}^n x_i}_{\hat{\mu}} \hat{\mu} + \underbrace{\frac{1}{n} \sum_{i=1}^n \hat{\mu} \cdot \hat{\mu}}_{\hat{\mu}^2}\right) \\
 &= E\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - 2 \hat{\mu}^2 + \hat{\mu}^2\right) \\
 &= E\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \hat{\mu}^2\right) \\
 &= E\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \mu^2 - \hat{\mu}^2 + \mu^2\right) \\
 &= E\left[\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \mu^2\right) - (\mu^2 - \hat{\mu}^2)\right] \\
 &= \frac{E\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \mu^2\right)}{A} - \frac{E(\hat{\mu}^2 - \mu^2)}{B} \quad \text{****}
 \end{aligned}$$

$$A: E\left(\frac{1}{n} \sum_{i=1}^n x_i^2 - \mu^2\right)$$

$$\begin{aligned}
 &= \frac{1}{n} \sum_{i=1}^n E(X_i^2) - E(\mu)^2 \\
 &= \frac{1}{n} \cdot n \cdot \sigma^2 \\
 &= \sigma^2
 \end{aligned}$$

$$B: E(\hat{\mu}^2 - \mu^2)$$

$$\begin{aligned}
 &= E(\hat{\mu}^2) - \mu^2 \\
 &= E(\hat{\mu}^2) - E(\hat{\mu})^2 \\
 &= \text{Var}(\hat{\mu})
 \end{aligned}$$

$E(\hat{\mu}) = \mu$
 but not $\hat{\mu} \neq \mu$
 not $E(\mu^2) = \mu^2$

$$\text{Var}(\hat{\mu}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right)$$

$$\begin{aligned}
 &= \frac{1}{n^2} n \cdot \sigma^2 \\
 &= \frac{\sigma^2}{n}
 \end{aligned}$$

$$\text{***} \Rightarrow \sigma^2 + \frac{\sigma^2}{n} = \left(\frac{n+1}{n}\right) \sigma^2 \neq \sigma^2 \Rightarrow \text{biased estimator.}$$

To correct:

$$\boxed{\frac{N}{N+1} \sigma^2}$$

→ unbiased estimator!

Useful properties of Gaussian

Assume $X \sim N(\mu_x, \sigma_x^2)$ ① Gaussian fits $\{x_1, x_2, \dots, x_N\}$ are Gaussian

$Y \sim N(\mu_y, \sigma_y^2)$ ② \vec{x} Gaussian \vec{y} also Gaussian

then.

③ $\vec{z} = \vec{a}\vec{x} + \vec{b}$ Gaussian

① Let $Z = aX + b$, then $Z \sim N(a\mu_x + b, a^2\sigma_x^2)$ (Linear transformer)

② Let $S = X + Y$, then $S \sim N(\mu_x + \mu_y, \sigma_x^2 + \sigma_y^2)$ (Sum and product)

③ Products of Gaussian r.v. are NOT Gaussian.

However, the product of Gaussian PDFs has Gaussian form.

$$\text{Norm PDF}(x | \mu_x, \sigma_x^2) = \text{Norm PDF}(x | \mu_y, \sigma_y^2) \perp \text{Norm PDF}(x | \mu, \sigma^2)$$

$$\text{where } \mu = \frac{\sigma_x^2 \mu_x + \sigma_y^2 \mu_y}{\sigma_x^2 + \sigma_y^2} \quad \sigma^2 = \frac{1}{\frac{1}{\sigma_x^2} + \frac{1}{\sigma_y^2}}$$

Covariance and covariance Matrices

$$\text{Cov}(X, Y) = E[(X - E(X))(Y - E(Y))] = E(XY) - E(X)E(Y)$$

If X is d dimensional random vector, its covariance matrix is
a [symmetric positive definite matrix]

$$\begin{aligned}\text{Cov}(X) &= E[(X - E(X))(X - E(X))^T] \\ &= \begin{pmatrix} \text{Var}(X_1) & \text{Cov}(X_1, X_2) & & \\ \text{Cov}(X_2, X_1) & \text{Var}(X_2) & & \\ & \ddots & \text{Var}(X_{d-1}, X_d) & \\ & & \text{Cov}(X_d, X_{d-1}) & \text{Var}(X_d) \end{pmatrix}\end{aligned}$$

Review: Positive Definite matrices (Gilbert P45)

Some properties:

- ① $\lambda_i > 0$
- ② $x^T S x$ is positive for all vector $x \neq 0$

Ex:

$$\text{Energy } x^T S x = \begin{bmatrix} x_1 & x_2 \end{bmatrix} \begin{bmatrix} 2 & 4 \\ 4 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 2x_1^2 + 8x_1x_2 + 9x_2^2$$

Is this positive for every x_1 and x_2 except $(x_1, x_2) = (0, 0)$? Yes, it is a sum of squares:

$$x^T S x = 2x_1^2 + 8x_1x_2 + 9x_2^2 = 2(x_1 + 2x_2)^2 + x_2^2 = \text{positive energy.}$$

③ $S = A^T A$ for matrix A with independent columns

④ All the leading determinants D_1, D_2, \dots, D_n of S are positive

CS336 Lecture 7

Bishop 2.3.1 - 2.3.6

Outline

1. Multivariate Gaussian Distribution
2. Covariance Properties: Why Contours are Elliptical
3. Margins of Gaussian are Gaussian
4. Conditionals of Gaussian are Gaussian
5. Linear-Gaussian models are some Gaussian.

Multivariate Gaussian Distribution

3.1

Joint Distribution of two independent gaussians

$$X_1 \sim N(\mu_1, \sigma_1^2) \quad \downarrow \quad \text{cov}(X_1, X_2) = 0$$

$$X_2 \sim N(\mu_2, \sigma_2^2)$$

$$\begin{aligned} p(X_1, X_2) &= \text{Norm PDF}(\mu_1, \sigma_1^2) \cdot \text{Norm PDF}(\mu_2, \sigma_2^2) \\ &= \frac{1}{\sqrt{2\pi}\sigma_1} e^{-\frac{(x-\mu_1)^2}{2\sigma_1^2}} \cdot \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{(x-\mu_2)^2}{2\sigma_2^2}} \\ &= \frac{1}{\sqrt{2\pi}\sigma_1} \cdot \frac{1}{\sqrt{2\pi}\sigma_2} \cdot e^{-\frac{1}{2}(x-\mu)^T \begin{bmatrix} \frac{1}{\sigma_1^2} & 0 \\ 0 & \frac{1}{\sigma_2^2} \end{bmatrix} (x-\mu)} \\ \text{where } X &= \begin{bmatrix} X_1 \\ X_2 \end{bmatrix} \quad \mu = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} \quad \Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \\ \Rightarrow &= \frac{1}{\sqrt{2\pi}\sigma_1} \cdot \frac{1}{\sqrt{2\pi}\sigma_2} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \\ &= \frac{1}{(2\pi)^{\frac{D}{2}}} \cdot \frac{1}{(2\pi)^{\frac{D}{2}}} \cdot \frac{1}{\sigma_1} \cdot \frac{1}{\sigma_2} \cdot e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)} \\ &= \frac{1}{(2\pi)^{D/2}} \cdot \frac{1}{|\Sigma|^{1/2}} \end{aligned}$$

Definition

Let X be a D -dim. random variable $X = [x_1, x_2, \dots, x_D]^T$

$$X \sim MVN(\mu, \Sigma)$$

Parameters:

= mean vector" $\mu \in \mathbb{R}^D$

= covariance" Σ is $D \times D$

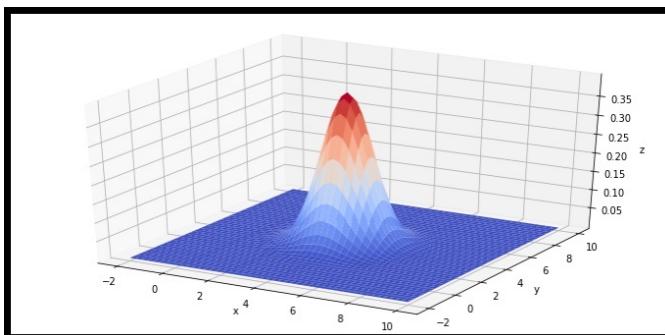
PDF:

$$MVN \text{ PDF } (X; \mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} \underline{(x-\mu)^T \Sigma^{-1} (x-\mu)} \right]$$

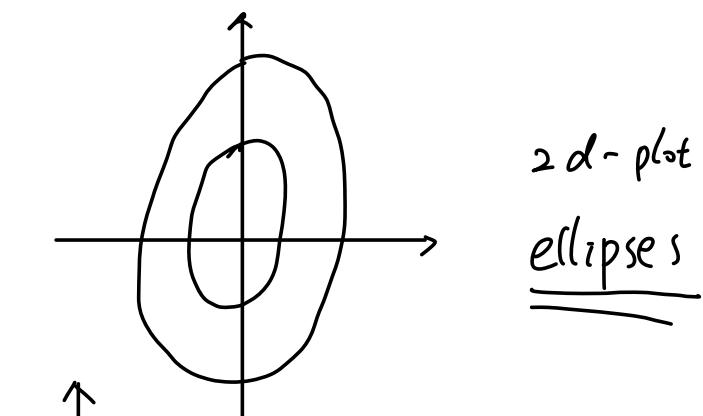
dimension check:

C : scalar \times scalar $\in \mathbb{R}$

$$f(x, \mu, \Sigma) : |xD \times D \times D \times D| = |x| \in \mathbb{R}$$



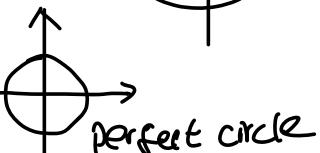
3d-plot



2d-plot

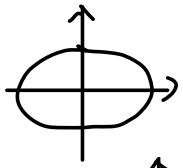
ellipses

$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

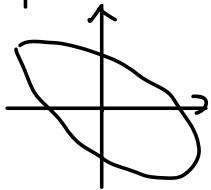


perfect circle

$$\Sigma = \begin{bmatrix} 3 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 3 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$



Some properties:

$$E(X) = \mu, \text{Var}(X) = \Sigma. \quad \boxed{E(X X^T) = \Sigma + \mu \mu^T}$$

Covariance Properties: Why Contours are Elliptical

(Symmetric Positive Definite Matrix)

由于方差矩阵 Σ 是正定且对称的，对其进行特征值分解：

$$\Sigma = Q \Lambda Q^{-1} = Q \Lambda Q^T, \text{ 其中 } QQ^T = Q^T Q = I, \Lambda = \text{diag}(\lambda_1, \lambda_2, \dots, \lambda_p)$$

$$\Rightarrow \Sigma = (q_1 \ q_2 \ \dots \ q_p) \cdot \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_p \end{bmatrix} \begin{bmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_p^T \end{bmatrix}$$

$$= \begin{bmatrix} q_1 \lambda_1 & q_2 \lambda_2 & \dots & q_p \lambda_p \end{bmatrix} \begin{bmatrix} q_1^T \\ q_2^T \\ \vdots \\ q_p^T \end{bmatrix} \rightarrow \Sigma$$

$$\Sigma = \sum_{i=1}^p q_i \lambda_i q_i^T$$

$$\Rightarrow \Sigma^{-1} = (Q \Lambda Q^T)^{-1}$$

$$= (Q^T)^{-1} \cdot \Lambda^{-1} \cdot Q^{-1} \quad \text{对正交矩阵 } Q^{-1} = Q^T$$

$$= Q \cdot \Lambda^{-1} \cdot Q^T \quad \Lambda^{-1} = \text{diag}(\frac{1}{\lambda_i})$$

$$\Sigma^{-1} = \sum_{i=1}^p q_i \cdot \frac{1}{\lambda_i} q_i^T$$

$$(X - \mu)^T \Sigma^{-1} (X - \mu) = (X - \mu)^T \sum_{i=1}^p q_i \cdot \frac{1}{\lambda_i} q_i^T (X - \mu)$$

$$= \sum_{i=1}^p \frac{1}{\lambda_i} (X - \mu)^T q_i \ q_i^T (X - \mu) *$$

$$\text{Let } (X - \mu)^T q_i = y_i. \text{ then } q_i^T (X - \mu) = y_i^T$$

$$* \Rightarrow \sum_{i=1}^p \frac{1}{\lambda_i} y_i y_i^T = \boxed{\sum_{i=1}^p \frac{y_i^2}{\lambda_i}}$$

$$\det D = 2 \Rightarrow \underbrace{\frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2}}_{=} = 1$$

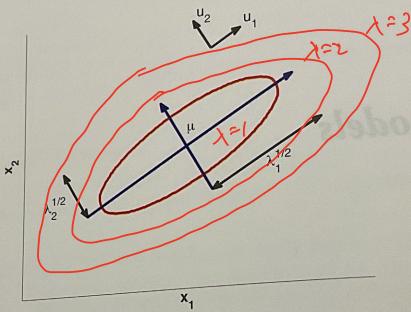


Figure 4.1 Visualization of a 2 dimensional Gaussian density. The major and minor axes of the ellipse are defined by the first two eigenvectors of the covariance matrix, namely \mathbf{u}_1 and \mathbf{u}_2 . Based on Figure 2.7 of (Bishop 2006).

The eigenvectors determine the orientation of the ellipse, and the eigenvalues determine how elongated it is.

Margins of Gaussian are Gaussian

Suppose we have D -dimensional cov. X

$$\mathbf{X} \sim \text{MVN}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \quad \boldsymbol{\mu} \in \mathbb{R}^D, \quad \boldsymbol{\Sigma} \in \mathbb{R}^{D \times D}$$

Consider any partition of indices

$$\{1, 2, 3, \dots, \underbrace{D-2, D-1, D}_{\text{first } D_A \text{ dims belong to A}}\} \quad \text{remaining dims belongs to B}$$

then $\begin{bmatrix} \mathbf{x}_A \\ \mathbf{x}_B \end{bmatrix} = \begin{bmatrix} \boldsymbol{\mu}_A \\ \boldsymbol{\mu}_B \end{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{AA} & \boldsymbol{\Sigma}_{AB} \\ \boldsymbol{\Sigma}_{BA} & \boldsymbol{\Sigma}_{BB} \end{bmatrix}$, where $\boldsymbol{\Sigma}_{AB} = \boldsymbol{\Sigma}_{BA}^T$

We have the margin of A:

$$\begin{aligned} p(X_A) &= \int p(X_A, X_B) d\mathbf{B} \\ &= \text{MVNorm PDF}(X_A | \boldsymbol{\mu}_A, \boldsymbol{\Sigma}_{AA}) \end{aligned}$$

By symmetry:

$$\begin{aligned} p(X_B) &= \int p(X_A, X_B) d\mathbf{A} \\ &= \text{MVNorm PDF}(X_B | \boldsymbol{\mu}_B, \boldsymbol{\Sigma}_{BB}) \end{aligned}$$

Conditionals of a Joint Gaussian are Gaussian

Given joint distribution over X_A, X_B

$$\begin{bmatrix} X_A \\ X_B \end{bmatrix} \sim MVN \left(\begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{AB}^T & \Sigma_{BB} \end{bmatrix} \right)$$

Define $\Lambda = \Sigma^{-1}$

$$\Rightarrow \Lambda = \begin{bmatrix} \Lambda_{AA} & \Lambda_{AB} \\ \Lambda_{AB}^T & \Lambda_{BB} \end{bmatrix}$$

Define the conditional density of X_A given $X_B = m_B$

$$P(X_A | X_B = m_B) = \frac{P(X_A, X_B = m_B)}{P(X_B = m_B)}$$

$$= \text{MultiNorm PDF} \left(X_A \mid \frac{\mu_A}{\Lambda_{AA}} - \frac{\Lambda_{AB}^T}{\Lambda_{AA}} \frac{\Lambda_{AB}(m_B - \mu_B)}{\Lambda_{BB}} + \frac{\Lambda_{AA}^{-1}}{\Lambda_{AA}} \right)$$

Linear Gaussian model is Gaussian

Consider a model with two variables :

$$1) \quad x \sim MVN(\mu, \Lambda^{-1}) \quad x \in \mathbb{R}^s$$

$$2) \quad y(x) \sim MVNL(Ax + b, L^{-1}) \quad y \in \mathbb{C}^{qT}$$

What is the joint distribution ? $P(x, y)$

$$P(x, y) = P(x) \cdot P(y|x)$$

$$\log P(x, y) = \log P(x) + \log P(y|x)$$

$$= \text{const} - \underbrace{\frac{1}{2}(x - \mu)^T \Lambda (x - \mu)}_A - \underbrace{\frac{1}{2}(y - Ax - b)^T L (y - Ax - b)}_B$$

$$A: \frac{1}{2}(x - \mu)^T \Lambda (x - \mu) = \frac{1}{2}(x^T \Lambda x - x^T \Lambda \mu - \mu^T \Lambda x + \frac{\mu^T \Lambda \mu}{c}) \\ = \underbrace{\frac{1}{2} x^T \Lambda x}_A - \underbrace{\frac{1}{2} \cdot 2(x^T \Lambda \mu)}_B$$

$$B: \frac{1}{2}(y - Ax - b)^T L (y - Ax - b) = \underbrace{-\frac{1}{2} x^T A^T L A x}_A - \underbrace{\frac{1}{2} (-2) y^T L A x}_B - \underbrace{\frac{1}{2} y^T L y}_C - \underbrace{\frac{1}{2} (2) x^T A^T L b}_D - \underbrace{\frac{1}{2} (-2) y^T L b}_E$$

order 2 order 1

$$\Rightarrow -\frac{1}{2} \underbrace{\begin{pmatrix} x \\ y \end{pmatrix}^T \begin{bmatrix} P \end{bmatrix} \begin{pmatrix} x \\ y \end{pmatrix}}_{\text{order } 2} + \underbrace{\begin{pmatrix} x \\ y \end{pmatrix}^T \begin{bmatrix} P_m \end{bmatrix}}_{\text{order } 1}$$

$$P = \begin{bmatrix} \Lambda + A^T L A & -A^T L \\ -L A & L \end{bmatrix}$$

$$m = P^{-1} \begin{bmatrix} \Lambda M - A^T L b \\ L b \end{bmatrix}$$

Lec 8

1. Probabilistic view of linear regression 3.1 - 3.3
2. ML estimation of weights and precision
3. Towards a full probabilistic model for regression
4. Posterior of Gaussian-Gaussian linear regression^M
5. MAP estimation of weights

Probabilistic view of linear regression

Goal: Given dataset $\{x_n, t_n\}_{n=1}^N$, with $x_n \in \mathbb{R}^D$ and $t_n \in \mathbb{R}$
want to predict t_* given x_*

Assume: "Linear model", which means prediction function is linear function
of input x_n

$$\begin{aligned}y(x_n, w) &= w_0 + w_1 x_{n1} + w_2 x_{n2} + \dots + w_D x_{nD} \\&= \sum_{d=0}^D w_d x_{nd} \quad (\text{define } x_{n0} = 1) \quad x_{nd} \in \mathbb{R}^{D+1} \\&= w^T x_n\end{aligned}$$

- often can define "smarter" features by transforming input x_n into another feature space via $\phi(x_n)$

$$\text{define } \phi(x_n) = [1 \quad \phi_1(n) \quad \phi_2(n) \quad \dots \quad \phi_M(n)]$$

Note: $\phi(x_n)$ can be nonlinear, e.g. Gaussian -- M total entries

- "Feature" model for prediction:

$$y(x_n, w) = \sum_{m=1}^M w_m \phi_m(x_n) = w^T \phi(x_n)$$

$$\boxed{\begin{aligned}w &= (w_0, \dots, w_{M-1})^T \\ \phi &= (\phi_0, \dots, \phi_{M-1})^T\end{aligned}}$$

Note: our prediction will not be perfect. Need to tolerate some noise. Let's define a probabilistic approach.

• likelihood of observing output t_n given input x_n $\epsilon \sim N(0, \sigma^2)$

$$p(t_n | x_n, w, \beta) = N(t_n | w^T \phi(x_n), \beta^{-1})$$

↑ mean variance $w^T \phi(x_n)$ is a scalar, linear
1-d variable in \mathbb{R} transformation of Gaussian is Gaussian
 $x_n \sim N(\mu, \sigma^2)$, $y = ax + b \sim N(ax + b, a^2 \sigma^2)$

If we assume all N observations are i.i.d. from this distribution, MLE:

$$p(t_n | x_n, w, \beta) = \prod_{n=1}^N (t_n | w^T \phi(x_n), \beta^{-1})$$

here $x \sim N(0, \sigma^2)$, $N(b, \frac{\sigma^2}{\beta^{-1}})$

$$\log p(t_n | x_n, w, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log (\pi_0) - \beta \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

key idea: use feature ϕ to get flexible prediction. Assume Gaussian noise model $N(t | w^T \phi(x_n), \beta^{-1})$

ML estimation of weights and precision

$$L(w, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log (\pi_0) - \beta \sum_{n=1}^N (t_n - w^T \phi(x_n))^2$$

$$\begin{aligned} \boxed{\frac{\partial L(w, \beta)}{\partial w}} &= -\frac{\beta}{2} \sum_{n=1}^N (t_n - w^T \phi(x_n))^2 \\ &\stackrel{\text{set}}{=} \mathbb{E}_{\text{mix}} \sum_{n=1}^N (t_n - w^T \phi(x_n)) \cdot \frac{\phi(x_n)^T}{\text{mix}} \stackrel{\text{set}}{=} 0 \\ &= \sum_{n=1}^N t_n \phi(x_n)^T - \sum_{n=1}^N w^T \phi(x_n) \phi(x_n)^T \stackrel{\text{set}}{=} 0 \\ &\Rightarrow \sum_{n=1}^N t_n \phi(x_n)^T = \sum_{n=1}^N w^T \phi(x_n) \phi(x_n)^T \quad * \\ &\text{Let } \Phi = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_N) \end{bmatrix}_{N \times M}, t = \begin{bmatrix} t_1 \\ \vdots \\ t_N \end{bmatrix}_{N \times 1} \end{aligned}$$

$$* \Rightarrow t^T \Phi = w^T (\Phi^T \Phi)$$

Transpose: $t^T \Phi^T = \Phi^T \Phi w$ $\Phi^T \Phi$ has to be invertible!

$$w = \frac{(\Phi^T \Phi)^{-1} \Phi^T t}{M \times N \quad N \times M \quad M \times M \quad N \times 1} \rightarrow \underline{M \times 1} \quad \checkmark$$

也可以用MLE方法！

$$J(w) = \|wx - y\| = \sum_{i=1}^n \|wx_i - y_i\|^2$$

- Penalized MLE

$$\max_w L(w, \beta) + \lambda \sum_{m=1}^M w_m^2 \quad (\text{Ridge Regression})$$

$$w^* = (\lambda I_M + \Xi^T \Xi)^{-1} \Xi^T t$$

this is always rank M and always invertible!

towards a full probabilistic model for regression

All unknown parameters (weight vector $w \in \mathbb{R}^M$ - precision β^{-1}) are treated probabilistically. For now, we'll assume β^{-1} is fixed known.

Equivalently, we need to define a joint model

$$P(t, w | X, \beta) \quad (\text{posterior} \propto \text{likelihood} \times \text{prior})$$

$$= P(t | w, X, \beta) \cdot P(w)$$

why? Given this joint, we can talk about posterior about parameters after seeing data:

$$P(w | \{x_n, t_n\}_{n=1}^N, \beta)$$

- can use the MAP estimate instead of ML
- can use samples from posterior to access uncertainty

We can also use posterior to make predictions about new data

use predictive posterior

$$P(t_* | x_*, \{x_n, t_n\}_{n=1}^N, \beta) = \int_w P(t_* | w, x_*, \{x_n, t_n\}_{n=1}^N, \beta) dw$$

Remark:

$$P(x=x_* | D) = \int_u P(x=x_*, u | D) du = \int_w P(t_* | w, x_*, \beta) \cdot P(w | \{x_n, t_n\}_{n=1}^N, \beta) dw$$

$$= \int_u P(x=x_* | u, D) \cdot P(u | D) du$$

$$\text{i.e.} = \int_u P(x=x_* | u) \cdot P(u | D) du$$

key difference: Average overall w vectors, weighted by posterior density. Do not just commit 100% to one w vector.

Review Gaussian

Given

1. Two independent Gaussian r.v.s

$$x \sim N(\mu_x, \sigma_x^2)$$

$$y \sim N(\mu_y, \sigma_y^2)$$

Their joint is Gaussian

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N\left[\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \sigma_x^2 & 0 \\ 0 & \sigma_y^2 \end{bmatrix}\right]$$

2. Two linearly-dependent Gaussian r.v.s

$$x \sim N(\mu_x, \Lambda^{-1})$$

$$y \sim N(\underline{\mu_x + b}, \Lambda^{-1})$$

Their joint is Gaussian

$$\begin{bmatrix} x \\ y \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix}, \begin{bmatrix} \Lambda^{-1} & -\Lambda^{-1} \\ -\Lambda^{-1} & \Lambda^{-1} \end{bmatrix}\right)$$

3. A joint distribution over a partitioned vector

$$\begin{bmatrix} \cdots & x & \cdots \end{bmatrix}^\top = \begin{bmatrix} \cdots & x_A & \cdots & x_B & \cdots \end{bmatrix}^\top$$

$$\begin{bmatrix} x_A \\ x_B \end{bmatrix} \sim N\left(\begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \Sigma_{AA} & \Sigma_{AB} \\ \Sigma_{BA} & \Sigma_{BB} \end{bmatrix}\right)$$

Marginal is Gaussian

$$p(x_A) = N(\mu_A, \Sigma_{AA})$$

Conditional is Gaussian

$$p(x_A | x_B) = N(\mu_{A|B}, \Sigma_{A|B})$$

$$\mu_{A|B} = \mu_A - \Sigma_{AA}^{-1} \Sigma_{AB} (x_B - \mu_B)$$

4. Two linearly-dependent Gaussians.

$$x \sim N(\mu, \Lambda^{-1})$$

$$y|x \sim N(Ax + b, L^{-1})$$

Posterior is Gaussian

$$p(x|y) = N(\mu_{xy}, \Sigma)$$

Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(x) = N(x|\mu, \Lambda^{-1}) \quad (2.113)$$

$$p(y|x) = N(y|Ax + b, L^{-1}) \quad (2.114)$$

the marginal distribution of y and the conditional distribution of x given y are given by

$$p(y) = N(y|A\mu + b, L^{-1} + A\Lambda^{-1}A^T) \quad (2.115)$$

$$p(x|y) = N(x|\Sigma\{A^T L(y - b) + \Lambda\mu\}, \Sigma) \quad (2.116)$$

where

$$\Sigma = (\Lambda + A^T L A)^{-1}. \quad (2.117)$$

Posterior of Gaussian - Gaussian (linear regression)

Assume $\beta^{-1} > 0$ known

Prior: $P(w) = N(m_0, S_0)$

$m_0 \in \mathbb{R}^m$

S_0 is $m \times m$ covariance matrix

p_{LX}

Likelihood:

$$p(t|x, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

$$p_{LY}(x) \quad \text{why? } t = w^T \phi(x) + \epsilon \rightarrow N(0, \beta^{-1})$$

$$\Rightarrow p(t|x, w, \beta^{-1}) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

Posterior:

$$p(w|t, x, \beta) = N(M_N, S_N)$$

$$p_{LX|Y} \quad \text{where } M_N = S_0^{-1} \underbrace{S_0^{-1} m_0 + \beta \bar{\Phi}^T t}_{\text{Gaussian} \times \text{Gaussian} \rightarrow \text{Gaussian}} \quad (3.50 - 3.51)$$

$$S_N^{-1} = S_0^{-1} + \beta \bar{\Phi}^T \bar{\Phi}$$

$$\text{where: } \bar{\Phi} = \begin{bmatrix} \phi(x_1) \\ \phi(x_2) \\ \vdots \\ \phi(x_n) \end{bmatrix}_{n \times m}, \quad t = \begin{bmatrix} t_1 \\ \vdots \\ t_n \end{bmatrix}_{n \times 1}$$

What is the MAP estimator for linear regression?

$$w_{MAP} = \underset{w \in \mathbb{R}^m}{\operatorname{argmax}} \log p(w | \{x_n, t_n\}_{n=1}^N) \leftarrow \text{we've shown this is } N(M_N, S_N)$$

$$\text{thus. } w_{MAP} = M_N$$

$$= \boxed{(S_0^{-1} + \beta \bar{\Phi}^T \bar{\Phi})^{-1} (S_0^{-1} m_0 + \beta \bar{\Phi}^T t)} *$$

Consider a prior that favors zero mean and diag covariance

$$\Rightarrow m_0 = \vec{0}, \quad S_0 = \begin{bmatrix} \sigma^2 & & & \\ & \ddots & & \\ & & \ddots & \\ & & & \sigma^2 \end{bmatrix} = \sigma^2 I_m, \quad S_0^{-1} = \frac{1}{\sigma^2} I \quad **$$

plug ** into * :

$$(\frac{1}{\sigma^2} I_m + \beta \bar{\Phi}^T \bar{\Phi})^{-1} \beta \bar{\Phi}^T t$$

$$= \frac{\beta}{\sigma^2} (\frac{1}{\beta} I_m + \bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T t$$

$$= \boxed{(\frac{\beta}{\sigma^2} I_m + \bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T t}$$

$$w^* = (\lambda I_m + \bar{\Phi}^T \bar{\Phi})^{-1} \bar{\Phi}^T t$$

Ridge Regression!

Equivalent to our sum-of-squares penalized linear regression

Conclusion:

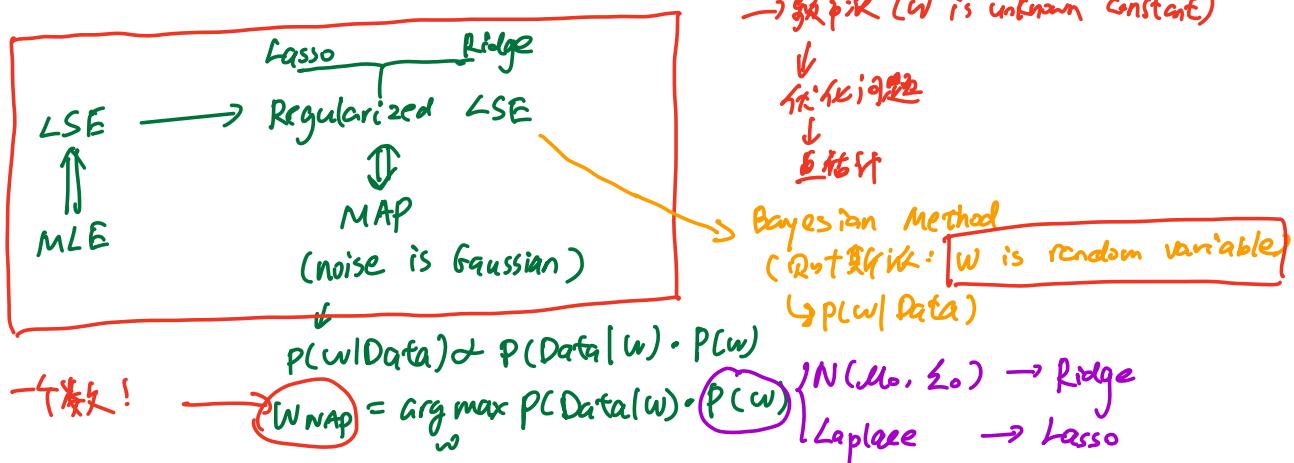
Our MAP for Gaussian-Gaussian model is equal to ridge regression,

under the constraint that

- prior is symmetric across w with precision λ
- known likelihood pre β

Supplement

1.



2. 算出来 $p(w|Data)$ 以后怎么求 T 值?

算 x^*, y^* ?

Inference: $p(w|Data)$
 $w|Data \sim N(\mu_w, \Sigma_w)$

$$\mu_w = \sigma^{-2} A^\top X^\top Y$$

$$\Sigma_w = A^{-1}$$

$$A = \sigma^{-2} X^\top X + \Sigma_p^{-1}$$

Prediction:

Given x^*, y^* . noise $\rightarrow x^* w \sim N(x^* \mu_w, x^* \Sigma_w x^*)$

Model: $f(x) = w^\top x = x^\top w$

$$\begin{cases} y = f(x) + \varepsilon \\ \varepsilon \sim N(0, \sigma^2) \end{cases}$$

Bayesian Linear Regression \rightarrow Prediction
 $p(f(x^*)) | Data, x^* = N(x^* \mu_w, x^* \Sigma_w x^*)$

$f(x) = x^\top w$

$f(x^*) = x^{*\top} w$

$p(w) = N(\mu_w, \Sigma_w)$

$w \sim N(\mu_w, \Sigma_w)$

$y^* = f(x^*) + \varepsilon$

$\varepsilon \sim N(0, \sigma^2)$

$\therefore p(y^* | Data, x^*) = N(x^{*\top} \mu_w, x^{*\top} \Sigma_w x^* + \sigma^2)$

shuhuai008 bili bili

3. 总结

Bayesian Linear Regression
shuhuai008 bilibili

Data: $\{(x_i, y_i)\}_{i=1}^n$, $x_i \in \mathbb{R}^p$, $y_i \in \mathbb{R}$.

Model: $\begin{cases} f(x) = w^T x = x^T w \\ y = f(x) + \epsilon \\ \epsilon \sim N(0, \sigma^2) \end{cases}$

Bayesian Method:

参数: w 不是未知的常量
 w 是一个概率分布

① Inference: $P(w | \text{Data}) \rightarrow \text{posterior}$

$P(w | \text{Data}) \propto \underbrace{P(w | \text{Data})}_{N(\mu_w, \Sigma_w)} \times \underbrace{\text{likelihood}}_{N(\Delta, \Delta)} \times \underbrace{\text{prior}}_{N(a, b)}$

$\mu_w = ?, \Sigma_w = ?$

② Prediction: Given x^*, y^* ?

$$P(y^* | \text{Data}, x^*) = \int_w \underbrace{P(y^* | w, \text{Data}, x^*)}_{P(y^* | w, x^*)} \cdot \underbrace{P(w | \text{Data})}_{\substack{\text{posterior}}} dw$$

Lec 9

1. Posterior predictive 3.3.1 - 3.3.2
2. Evidence = A measure of model quality 3.5
3. Model selection 3.4
4. Hyperparameter estimation

Posterior predictive

Recall our model:

Prior: $p(w|\alpha) = N(w|0, \alpha^{-1}I)$ zero-mean isotropic Gaussian

Likelihood: $p(t|w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$

iid assumption, says if we saw a new observation at x_* , then

$$p(t_*, t | w, \beta) = p(t | w, \beta) \cdot p(t_* | w, \beta)$$

Then by linear Gaussian rules:

Joint distribution $p(t_{1:N}, w | \alpha, \beta)$ is Gaussian

$$p(w, t) = \frac{p(w)}{\Gamma} \cdot \frac{p(t|w)}{\Gamma}$$

Expanded Joint distribution $p(t_*, t_{1:N}, w | \alpha, \beta)$ is Gaussian *

Marginal of * : $p(t_*, t_{1:N} | \alpha, \beta)$ is Gaussian. **

Condition of ** : $p(t_* | t_{1:N}, \alpha, \beta)$ is Gaussian

predictive posterior = $p(\text{new data} | \text{train data, hyperparameters})$

Formulas for Predictive Posterior

$$p(t_* | t_{1:N}, \alpha, \beta) = N(t_* | M_N^\top \phi(X_*)^T, \frac{1}{\beta} + \phi(X_*)^T S_N \phi(X_*))$$

where M_N M_{K1} is posterior mean vector

S_N $M_{K1} \times M_{K1}$ is posterior covariance vector

$\phi(X_*)$ $M \times 1$ is feature vector at test point X_*

β $|X|$ is scalar (like likelihood precision)

Useful Properties

- Average over many weights w , do not commit to a point estimate
 $p(t_* | t) = \int (t_*, w | t) dw = \int p(t_* | w, t) \cdot p(w | t) dw$
 $= \int p(t_* | w) \cdot p(w | t) dw$ (iid)
- Variance may change with location of prediction X_* . Will always be at least $\frac{1}{\beta}$. Can be larger. $\frac{1}{\beta} + \phi(X_*)^T S_N \phi(X_*) \geq 0$
- Variance $\sigma_N^2(X_*)$ cannot increase as more data seen.

Evidence : A measure of model quality

$\star w \in \mathbb{R}^M$

t all N observations. $t_n \in \mathbb{R}$

fixed hyperparameters

$\alpha > 0$

$\beta > 0$

free to add in !!!



- $p(w)$ prior $p(w | \alpha, \beta)$
- $p(t | w)$ likelihood $p(t | w, \alpha)$
- $p(w | t)$ posterior $p(w | t, \alpha, \beta)$
- $p(t)$ evidence $p(t | \alpha, \beta)$

Intuition: Similar to posterior predictive except we weight by prior instead of posterior. We do it for the whole given dataset.

$$p(t|\alpha, \beta) = \int \underbrace{p(t|w, \beta)}_{\text{likelihood}} \cdot \underbrace{p(w|\alpha)}_{\text{prior}} dw \quad \hookrightarrow$$

where: prior $p(w|\alpha) = N(w|0, \alpha^{-1} I)$ $p(y) = \int p(y|x) \cdot p(x) dx$
 likelihood $p(t|w, \beta) = N(t_n | w^T \phi(x_n), \beta^{-1})$ marginal probability

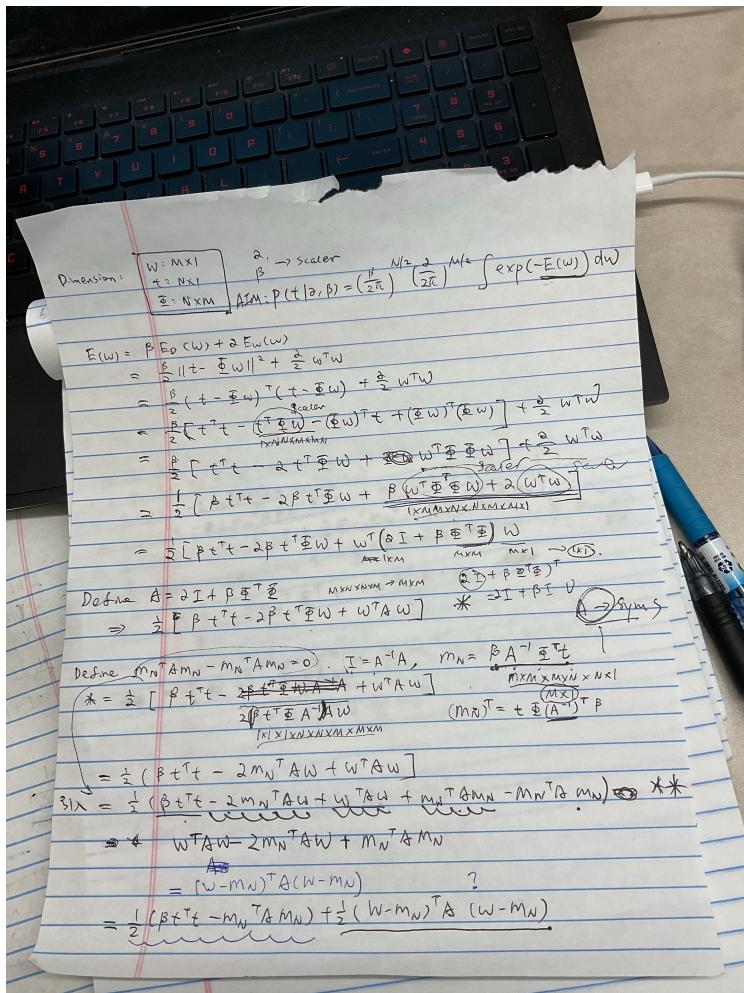
$$\text{Predictive: } p(t_* | t_{1:N})$$

$$= p(t_* | t_{1:N}, X_*)$$

$$= \int_w \underbrace{p(t_* | w, X_*)}_{\text{likelihood}} \cdot \underbrace{p(w | \{X_n, t_n\}_{n=1}^N)}_{\text{Posterior}} dw$$

结束

$$\ln p(\mathbf{t} | \alpha, \beta) = \frac{M}{2} \ln \alpha + \frac{N}{2} \ln \beta - E(m_N) - \frac{1}{2} \ln |\mathbf{A}| - \frac{N}{2} \ln(2\pi)$$



hint:

$$p(X) = N(\mu, \Sigma)$$

高斯分布:

$$\begin{aligned} & \exp(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)) \\ &= -\frac{1}{2}(x^T \Sigma^{-1} - \mu^T \Sigma^{-1})(x-\mu) \\ &= -\frac{1}{2}(x^T \Sigma^{-1} x - x^T \Sigma^{-1} \mu - \mu^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu) \\ &= -\frac{1}{2}(x^T \Sigma^{-1} x - 2 \mu^T \Sigma^{-1} x + \mu^T \Sigma^{-1} \mu) \\ &\quad \downarrow \qquad \downarrow \\ &= -\frac{1}{2} x^T \Sigma^{-1} x - \frac{1}{2} \mu^T \Sigma^{-1} \mu \end{aligned}$$

$L(X^k = n | Y = m) = \dots X^k = X^m = \frac{M^k M^m}{M^{k+m}}$
 $b(X^k = n | Y = m) = b(X^k | Y = m) = \frac{b(X^k)}{b(Y = m)}$
 $b(X^k + X^m) = b(X^k | Y = m) b(X^m | Y = m)$

$A = \alpha I + \beta \bar{A}^T \bar{A}$
 $\text{Define: } M_N^T A M_N - M_N^T \beta M_N = 0 \quad \Rightarrow \quad I = A^{-1} A$
 $M_N = \beta A^{-1} \bar{A}^T \bar{A}$

~~推上:~~
 $\left[\frac{1}{2} (\beta \bar{A}^T \bar{A} - M_N^T A M_N) + \frac{1}{2} (W - M_N)^T A (W - M_N) \right]$
 $\frac{1}{2} (\beta \bar{A}^T \bar{A} - M_N^T A M_N) \xrightarrow{\text{BTW}} \frac{1}{2} \| \bar{A}^T \bar{A} - M_N \|_F^2 + \frac{1}{2} M_N^T M_N \quad \text{没有 } A$
 $= \frac{1}{2} (\beta \bar{A}^T \bar{A} - M_N^T A M_N + M_N^T A M_N - M_N^T A M_N)$
 $= \frac{1}{2} (\beta \bar{A}^T \bar{A} - 2 M_N^T A M_N + M_N^T A M_N)$
 $\leq \frac{1}{2} (\beta \bar{A}^T \bar{A} - 2 M_N^T A \beta A^{-1} \bar{A}^T \bar{A} + M_N^T (\alpha I + \beta \bar{A}^T \bar{A}) M_N)$
 $= \frac{1}{2} (\beta \bar{A}^T \bar{A} - 2 M_N^T \bar{A}^T \bar{A} + 2 M_N^T M_N + \beta M_N^T \bar{A}^T \bar{A} M_N)$
 $= \frac{1}{2} (\cancel{\beta \bar{A}^T \bar{A}} - 2 M_N^T \bar{A}^T \bar{A} + \cancel{\beta M_N^T \bar{A}^T \bar{A} M_N}) - \frac{1}{2} M_N^T M_N$
 $\Rightarrow \bar{A}^T \bar{A} - 2 M_N^T \bar{A}^T \bar{A} + M_N^T \bar{A}^T \bar{A} M_N \quad \cancel{-}$
 $\Rightarrow \alpha^2 = 2ab \quad b^2$
 $\text{where } a = \bar{A}^T \bar{A} \quad b = M_N^T \bar{A}^T \bar{A}$
 $\Rightarrow \frac{\beta}{2} \| \bar{A}^T \bar{A} - M_N \|_F^2 + \frac{1}{2} M_N^T M_N \quad \text{proved}$

$E(W) = \frac{\beta}{2} \| \bar{A}^T \bar{A} - M_N \|_F^2 + \frac{1}{2} M_N^T M_N \quad \text{define as}$

~~推上:~~
 $\int \exp(-E(W)) dW$
 $= \int \exp\left(-\frac{\beta}{2} \| \bar{A}^T \bar{A} - M_N \|_F^2 + \frac{1}{2} M_N^T M_N\right)$

Given $S_n^{-1} = \alpha I + \beta \bar{A}^T \bar{A}$, means $A = S_n^{-1}$
 $\Rightarrow M_N = \beta S_n^{-1} \bar{A}^T \bar{A} \Leftarrow \text{高斯分布的均值}$

推定值有部分待证。

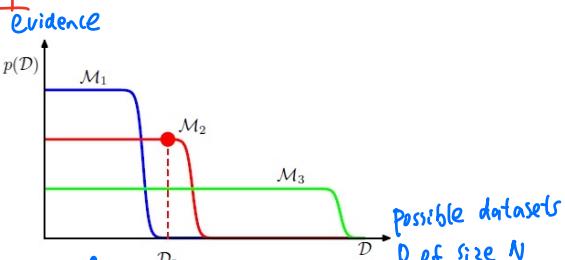
Model Selection

(Fixed validation/Cross Validation / Evidence)

use the evidence for selection

have unit area ←

Figure 3.13 Schematic illustration of the distribution of data sets for three models of different complexity, in which M_1 is the simplest and M_3 is the most complex. Note that the distributions are normalized. In this example, for the particular observed data set D_0 , the model M_2 with intermediate complexity has the largest evidence.



- simpler models give higher density to dataset they favor.
- Using evidence will avoid overfitting

Usual Procedure: Pick single best model M^* for L possible. Use that to make predictions.

Bayesian Procedure: Average all models, weighted by posterior:

$$p(t^*|t) = \sum_{\ell=1}^L p(t^*|t, M_\ell) \cdot p(M_\ell|t)$$

↑ predictive posterior ↑ model posterior
 for model ℓ

Comparison of Methods for Hyperparameter Selection

	Fixed valid. set (fraction f)	K-fold cross-validation	Bayesian evidence
Fraction data used for training run	$(1-f)$	$(K-1)/K$	100% 1.0 ✓ Higher is better Better use of training data
Total runs/examples seen for training	1 run $(1-f)N$ ✓	K runs $(K-1)*N$	1 run N ✓ Lower is better Faster training
Total runs/examples seen for evaluation of fitness	1 run fN	K runs N	1 run N ✓ Lower is better Faster evaluation
Fitness function	Heldout likelihood <i>Tow!</i>	Heldout likelihood <i>low</i>	Evidence <i>HARD</i>

Hyperparameter estimation

Want to know good values for α, β given dataset.

$$\hat{\alpha}, \hat{\beta} = \arg \max_{\alpha > 0, \beta > 0} \log P(\text{data} | \alpha, \beta)$$

This is "empirical Bayesian" point estimation

Can be solved via:

- ① enumeration via grid search
- ② gradient descent
- ③ coordinate descent (will see in EM algorithm)
- ④ analytical estimates

Lec 10 Probabilistic Generalized Linear models

1. Discriminative & Generative model
2. Generalized Linear models
3. Sigmoid function
4. Probabilistic Logistic Regression: ML & MAP strategy
5. 2nd-order gradient methods for linear + logistic Regression

Discriminative & Generative model

Consider Supervised Learning task:

Given training data $\{x_n, y_n\}_{n=1}^N$. Learn to predict $P(y_* | x_*)$

Two approaches to probabilistic modeling.

Discriminative (判别模型)

Model:
$$P(y|x) = \prod_{n=1}^N P(y_n|x_n)$$

(x is treated as fixed / known)

Parameter:

w generates model given feature

prediction: $P(y_* | x_*, w)$

(directly use likelihood)

Training: $\max_w \prod_n P(y_n | x_n, w)$

Pros: - simpler, fewer parameters

- directly solve supervised task

Cons: - cannot predict if x has missing values.

Generative (生成模型)

Model:
$$P(x, y) = \prod_{i=1}^N P_a(x_i | y_i) \cdot P_\pi(y_i)$$

(both x & y are r.v.)

Parameters:

Q : generate x given y

π : generate y

Prediction via Bayes rule:

$$\begin{aligned} P(y_* | x_*) &= \frac{P(x_* | y_*) \cdot P(y_*)}{P(x_*)} \\ &= \frac{P(x_* | y_*) \cdot P(y_*)}{P(x_*(y')) \cdot P(y')} \end{aligned}$$

Training: $\arg \max_{Q, \pi} \log P_a(x_n | y_n) + \log P_\pi(y_n)$

Pros: - can predict if x/w is missing value
- can sample new x from $P(x|y)$

Cons: - more complex
- more parameters

Generalized Linear models

Extend linear regression to tasks with constrained output space

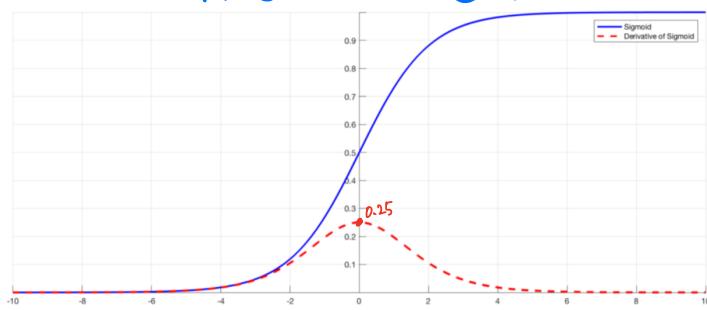
Training Data: $\{(x_n, t_n)\}_{n=1}^N$

Model: $p(t_{1:N} | X_{1:N}) = \prod_{n=1}^N Q(t_n | f(w^T \phi(x_n)))$

Output Space γ	name	distribution Q	link function f
real $(-\infty, +\infty)$	linear regression	Normal	$f(w, x) = w^T \phi(x)$
integers $\{0, 1, 2, \dots\}$	Poisson regression	Poisson	$f(w, x) = e^{w^T \phi(x)}$
binary $\{0, 1\}$	Logistic regression	Bernoulli	$f(w, x) = \frac{\sigma(w^T \phi(x))}{\text{logistic sigmoid}}$
	Probit regression	Bernoulli	$f(w, x) = \frac{\Phi(w^T \phi(x))}{\text{Normal CDF}}$
C-ary classification $\begin{matrix} 0 & 0 & 0 & \dots & 1 \\ 0 & 0 & 1 & \dots & 0 \\ 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & \dots & 0 \end{matrix}$	Multi-class logistic regression	Categorical	$f(w, x) = \text{Softmax}(W^T \phi(x))$ W: MxC matrix vector in R^C
C levels e.g. ratings 0-5 stars of item ranking e.g. cold, warm, hot	Ordinal regression	Ordinal	look up if interested!
positive reals $(0, +\infty)$	Exponential regression	Exponential	

Sigmoid Function

$$\sigma(r) = \frac{1}{1 + e^{-r}} = \frac{e^r}{e^r + 1}$$



Properties:

- maps reals to probability interval $(0, 1)$
- monotonic: $r_B > r_A \Rightarrow \sigma(r_B) > \sigma(r_A)$
- Invertible: reals $(-\infty, +\infty)$ $\xrightarrow[\text{logit}]{\text{Sigmoid}}$ probability $(0, 1)$ $\xrightarrow{\text{logit: } \ln \frac{p}{1-p}}$
- $1 - \sigma(r) = \sigma(-r)$

$$\sigma'(r) = \sigma(r)(1 - \sigma(r)) \rightarrow [0, 0.25]$$

Probabilistic Logistic Regression: ML & MAP strategy

$$\text{Prior: } p(w) = N(w | m_0, S_0) \\ = N(w | 0, \sigma^2 I_M) \text{ usually}$$

$$\text{Likelihood: } p(t_n | w) = \underbrace{\text{BernPMF}(t_n | \sigma(w^T \phi(x_n)))}_{\text{Bernoulli}}$$

Point Estimation Strategies:

$$\text{MLE: } \arg \min_{w \in \mathbb{R}^M} \sum_{i=1}^N \log p(t_i | w)$$

$$\text{MAP: } \arg \max_{w \in \mathbb{R}^M} \sum_{i=1}^N \log p(t_i | w) + \log p(w)$$

Next time: Posterior Estimation $p(w | t_{1:N})$

Strategies for optimization.

$\arg \min_{w \in \mathbb{R}^M} d(w)$ always $\arg \min - \arg \max$ to $\arg \min$

Strategies for Optimization

$$\arg \min_{w \in \mathbb{R}^M} d(w)$$

Convention: always minimize

Analytical methods

- Set up $\nabla_w d = 0$
- Manipulate to find $w^* = \dots$ closed form expression

ML/MAP for Linear Regression

ML/MAP for Logistic Reg.

No closed form solution exists.

Gradient methods

step size $\epsilon_t > 0$
gradient $g \in \mathbb{R}^M$
Hessian $H \in \mathbb{R}^{M \times M}$

- 1st order gradient descent
 $w_{t+1} \leftarrow w_t - \epsilon_t g(w_t)$

Many iterations,
each cheap:
 $O(M)$

- 2nd order gradient descent
 $w_{t+1} \leftarrow w_t - \epsilon_t H(w_t)^{-1} g(w_t)$

Few iter.,
each expensive:
 $O(M^3)$
inverting
matrix



Other methods

- grid search
- random search
- Nelder-Mead
- ... many others possible

usually very
inefficient in
high dimensions
($M > 3$)



2nd-order gradient methods for (linear + logistic) Regression

ML/MAP estimation: Gradients + Hessians

$$\text{For Logistic Regression: } \mathcal{L}(\omega) = - \sum_{n=1}^N \left[\log \text{BernPMF}(t_n | \sigma(w^T \Phi(x_n))) - \log \text{NormalPDF}(u | 0, \alpha^{-1} I_M) \right]$$

<p style="text-align: center;">gradient $\nabla_w \mathcal{L}$ shape: $(M, 1)$</p> <p>Linear Regr. $\beta \bar{\Phi}^T (\underbrace{\bar{\Phi} w - t}_{\substack{\text{predicted} \\ \text{label}}}) + \alpha w$</p> <p>Logistic Regr. $\bar{\Phi}^T \left(\underbrace{\sigma(\bar{\Phi} w)}_{\substack{\text{prediction} \\ \text{(probability)}}} - t \right) + \alpha w$</p>	<p style="text-align: center;">Hessian $\nabla_w \nabla_w \mathcal{L}$ shape: (M, M)</p> <p>$\beta \bar{\Phi}^T \bar{\Phi} + \alpha I_M$</p> <p>$\bar{\Phi}^T R(\omega) \bar{\Phi} + \alpha I_M$</p>
------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

$\mathcal{L}(\omega)$ has unique minimum w^* exists

Question:

Unique minima w^* exists when objective is convex

• Is Linear Regression MAP convex?

A: Yes. $\beta \bar{\Phi}^T \bar{\Phi} + \alpha I_M$
 P.S.D if $\beta > 0$

• Is LR MAP convex?

A: Yes. $\bar{\Phi}^T R(\omega) \bar{\Phi} + \alpha I_M$
 $R(\omega) \rightarrow$ positive diagonal
 $\bar{\Phi}^T R(\omega) \bar{\Phi}$ is still positive

2nd order gradient methods

solve an optimization Problem:

$$\underset{w \in \mathbb{R}^n}{\operatorname{argmin}} L(w)$$

$$\text{Let } g(w) = \nabla_w L \quad (M, 1)$$

$$H(w) = \nabla_w \nabla_w L \quad (M, M)$$

1st. order GD:

$$w^{\text{new}} = w^{\text{old}} - \varepsilon g(w^{\text{old}})$$

2nd order GD:

$$w^{\text{new}} = w^{\text{old}} - \varepsilon H(w^{\text{old}})^{-1} g(w^{\text{old}})$$

Linear Regression $\hat{\circ}$ 2nd order GD

Using formulas from previous page:

$$\begin{aligned} w^{\text{new}} &= w^{\text{old}} - \varepsilon \left[\beta \bar{\Phi}^T \bar{\Phi} + \alpha I_M \right]^{-1} \left(\beta \bar{\Phi}^T \bar{\Phi} w^{\text{old}} - \beta \bar{\Phi}^T t + \alpha w^{\text{old}} \right) \\ &= w^{\text{old}} - \varepsilon \left(\beta \bar{\Phi}^T \bar{\Phi} + \alpha I_M \right)^{-1} \left[\underbrace{\left(\beta \bar{\Phi}^T \bar{\Phi} + \alpha I_M \right) w^{\text{old}}}_{\text{MAP estimate}} - \beta \bar{\Phi}^T t \right] \\ &= w^{\text{old}} - \varepsilon w^{\text{old}} + \varepsilon \left(\beta \bar{\Phi}^T \bar{\Phi} + \alpha I_M \right)^{-1} \beta \bar{\Phi}^T t \end{aligned}$$

If step size $\varepsilon=1$, we get:

$$w^{\text{new}} = \left(\bar{\Phi}^T \bar{\Phi} + \frac{\lambda}{\beta} I_M \right)^{-1} \bar{\Phi}^T t \quad \begin{matrix} \text{optimal MAP} \\ \text{estimate} \\ \text{in one step!} \end{matrix}$$

Anytime loss is quadratic, 2nd order gradient update with step size $\varepsilon=1$ will find global minima in one step.

Logistic Regression w/ 2nd order GD

Still gold standard for ML/MAP estimation,
but requires many iterations (unlike linear regression ^{ML estimation using 2nd order methods}).

Update looks like this (assuming ML estimation. MAP is similar.)

$$w^{\text{new}} \leftarrow w^{\text{old}} - \varepsilon \left(\underline{\Phi}^T R(w^{\text{old}}) \underline{\Phi} \right)^{-1} \underline{\Phi}^T \left(\underbrace{r(\underline{\Phi} w)}_{\hat{y}} - t \right)$$

Assume $\varepsilon=1$ and rearranging terms

$$\begin{aligned} & \leftarrow \left(\underline{\Phi}^T R \underline{\Phi} \right)^{-1} \left(\underline{\Phi}^T R \underline{\Phi} w^{\text{old}} - \underline{\Phi}^T \hat{y} + \underline{\Phi}^T t \right) \\ & \leftarrow \left(\underline{\Phi}^T R \underline{\Phi} \right)^{-1} \left(\underline{\Phi}^T R z \right) \quad \text{where } z \in \mathbb{R}^N \\ & \qquad \qquad \qquad z = \underline{\Phi} w^{\text{old}} - R(\hat{y} - t) \end{aligned}$$

Interpret as "least squares"
with per-example "weights" r , $R = \text{diag}(r)$
and outcomes \underline{z} , both of which depend on
previous value of w .

Thus, the name

"Iteratively Re-weighted Least Squares" (IRLS)

Often used for Generalized Linear Models
to do ML estimation of weights $w \in \mathbb{R}^M$

while not converged:

update R given w
update \underline{z} given w
update $w \leftarrow (\underline{\Phi}^T R \underline{\Phi})^{-1} \underline{\Phi}^T R \underline{z}$

Supplement. 演習問題 (Bishop 207 - 208)

$$W_{\text{new}} = W^{\text{old}} - H^{-1} \triangledown E(w)$$

Let $\phi(X_n) = \phi_n$ Loss Function

$E(w)$: the sum-of-square error function
 Linear/Logistic Regression 與 線性回歸
 二元分類問題.

Linear Regression

$$E(w) = \frac{1}{2} \sum_{n=1}^N (t_n - w^T \phi_n)^2$$

$$\begin{aligned} \triangledown E(w) &= \sum_{n=1}^N (t_n - w^T \phi_n) \cdot (-\phi_n) \\ &= \sum_{n=1}^N (w^T \phi_n - t_n) \cdot (\phi_n) \\ &= w^T \sum_{n=1}^N \phi_n \phi_n^T - \sum_{n=1}^N t_n \phi_n \end{aligned}$$

$$\triangledown E(w) = \Phi^T \Phi w - \Phi^T t$$

$$H = \Phi \triangledown E(w) = \Phi^T \Phi$$

$$\Phi = \begin{bmatrix} -\varphi_1 - \\ -\varphi_2 - \\ \vdots \\ -\varphi_n - \end{bmatrix}_{m \times n} \quad w = \begin{bmatrix} \quad \\ \quad \end{bmatrix}_{n \times 1}$$

$$\begin{aligned} w_{\text{new}} &= W^{\text{old}} - H^{-1} \triangledown E(w) \\ &= W^{\text{old}} - (\Phi^T \Phi)^{-1} (\Phi^T \Phi w^{\text{old}} - \Phi^T t) \\ &= W^{\text{old}} - [W^{\text{old}} - (\Phi^T \Phi)^{-1} \Phi^T t] \\ &= (\Phi^T \Phi)^{-1} \Phi^T t \end{aligned}$$

Logistic Regression

$$w_{\text{new}} = w^{\text{old}} - H^{-1} \triangledown E(w)$$

Bernoulli

$$p(t|w) = \prod_{n=1}^N y_n^{t_n} (1-y_n)^{1-t_n} \quad t_n \in \{0, 1\}$$

$$E(w) = -\ln p(t|w) = -\sum_{n=1}^N (t_n \ln y_n + (1-t_n) \ln (1-y_n)) \quad *$$

$$\text{where } y_n = \sigma(w^T \phi_n)$$

$$* = -\sum_{n=1}^N \left(\underbrace{t_n \ln \sigma(w^T \phi_n)}_{①} + \underbrace{(1-t_n) \ln (1-\sigma(w^T \phi_n))}_{②} \right) \quad ***$$

$$\text{Grafen } \delta'(x) = \delta(1-\delta)$$

$$\begin{aligned}
\textcircled{1} \quad & t_n \nabla_{\omega} \delta(\omega^T \phi_n) \\
&= t_n \cdot \frac{1}{\sigma(\omega^T \phi_n)} \cdot \sigma(\cancel{\omega^T \phi_n}) \circ (1 - \sigma(\omega^T \phi_n)) \cdot \phi_n \\
&= \sum_{i=1}^n t_n \phi_n (1 - \delta(\omega^T \phi_n)) = \sum_{i=1}^n t_n \phi_n - t_n \phi_n \delta(\omega^T \phi_n) \quad \textcircled{1} \\
\textcircled{2} \quad & (1-t_n) \nabla_{\omega} (1 - \delta(\omega^T \phi_n)) \\
&= (1-t_n) \cdot \frac{1}{1 - \sigma(\omega^T \phi_n)} \cdot -(\sigma(\omega^T \phi_n) \cdot \cancel{(1 - \sigma(\omega^T \phi_n))} \cdot \phi_n) \\
&= \sum_{i=1}^n -((1-t_n) \cdot \phi_n \circ \sigma(\omega^T \phi_n)) \\
&= \sum_{i=1}^n (\phi_n t_n - \phi_n) \cdot \sigma(\omega^T \phi_n) \\
&= \sum_{i=1}^n \phi_n t_n \cancel{\delta(\omega^T \phi_n)} - \phi_n \delta(\omega^T \phi_n) \quad \textcircled{2}
\end{aligned}$$

$$\text{If } \textcircled{2}: - \sum_{i=1}^n (t_n \phi_n - \phi_n \delta(\omega^T \phi_n))$$

$$\begin{aligned}
\nabla E(\omega) &= \sum_{n=1}^n (y_n - t_n) \phi_n \\
&= \sum_{n=1}^n y_n \phi_n - \sum_{i=1}^n \phi_n t_n
\end{aligned}$$

$$\boxed{\nabla E(\omega) = \Phi^T (\mathbf{y} - \mathbf{t})}$$

$$\begin{aligned}
\nabla \nabla E(\omega) &= \sum_{n=1}^n (\phi_n \delta(\omega^T \phi_n) - t_n \phi_n) \\
&= \sum_{n=1}^n (\phi_n \sigma(\omega^T \phi_n) ((1 - \delta(\omega^T \phi_n)) \cdot \phi_n)) \\
&= \sum_{n=1}^n (\phi_n^T y (1-y) \cdot \phi_n) \\
&= \sum_{n=1}^n (\phi_n \phi_n^T) y (1-y) = \Phi^T R \Phi
\end{aligned}$$

$$\boxed{\nabla \mathcal{L}(w) = \Phi^T R \Phi}, \text{ where } R = \sum_n r_n (-y_n)$$

$$\begin{aligned} w_{\text{new}} &= w_{\text{old}} - H^{-1} \nabla \mathcal{L}(w) \\ &= w_{\text{old}} - (\Phi^T R \Phi)^{-1} (\Phi^T (y - t)) \\ &= w_{\text{old}} - (\Phi^T R \Phi)^{-1} (\Phi^T y - \Phi^T t) \end{aligned}$$

Update looks like this (assuming ML estimation. MAP is similar.)

$$w_{\text{new}} \leftarrow w_{\text{old}} - \varepsilon \left(\Phi^T R (w_{\text{old}}) \Phi \right)^{-1} \Phi^T \left(\underbrace{\tau(\Phi w)}_{\hat{y}} - t \right)$$

Assume $\varepsilon = 1$ and rearranging terms

$$\begin{aligned} &\leftarrow \left(\Phi^T R \Phi \right)^{-1} \left(\Phi^T R w_{\text{old}} - \Phi^T \hat{y} + \Phi^T t \right) \\ &\leftarrow \left(\Phi^T R \Phi \right)^{-1} \left(\Phi^T R z \right) \quad \text{where } z \in \mathbb{R}^N \\ &\quad z = \Phi w_{\text{old}} - R(\hat{y} - t) \end{aligned}$$

Interpret as "least squares" with per-example "weights" r , $R = \text{diag}(r)$

and outcomes \mathbf{z} , both of which depend on previous value of w .

Thus, the name

"Iteratively Re-weighted Least Squares" (IRLS)

Often used for Generalized Linear Models to do ML estimation of weights $w \in \mathbb{R}^M$

Lec 11 Bayesian Logistic Regression 4.4 . 4.5

1. Overview of posterior + predictive
 2. Laplace approximation in 1-dim and M-dim
 3. Laplace approximation for Logistic regression
 4. Posterior predictive approximation.

Overview of posterior + predictive

Model

prior on weights $\mathcal{W} \in \mathbb{R}^M$

$$p(w) = N(w | \mu_0, \Sigma_0)$$

Likelihood of "output" $t_n \in \{0, 1\}$

$$p(t|w) = \prod_{n=1}^N \text{Bern}(t_n | \sigma(w^T \phi(x_n))) \quad (\text{iid})$$

Goal: To estimate the Posterior and Predictive

Posterior: $p(w|t_{1:N})$

No closed form!

No closed form!
Not a Gaussian (Not Gau x Gau again!)

Predictive: $P(t_* | t_{1:N})$

$$= p(t_* | t_{1:N}, x_*)$$

$$= \int_w p(t_*|w, x_*) \cdot p(w| \{x_n, t_n\}_{n=1}^N) dw$$

w likelihood Posterior

- Must be a Bernoulli
 - But no close form for its parameter

Laplace approximation in 1-dim and M-dim

Given: a random variable $w \in \mathbb{R}$

$$p(w) = \frac{1}{Z} f(w)$$

Here, $f(w) > 0$ is known and differentiable
but computing $\int_w f(w) dw$ is hard

R: How can we estimate the distribution $p(w)$?

Intuition: $\ln p(w) \approx g(w) \approx \ln f(w) + p(w)$

Detail: • pick mean to match the mode of $p(w)$
 $m = \arg \max_w p(w) = \arg \max_w f(w)$

Can use gradient methods to solve this numerically

• pick precision $\frac{1}{\sigma^2}$ to perform best possible 2nd order
Taylor approximation to $p(w)$ at the mode $w=m$

$$\beta = \frac{\partial}{\partial w} \frac{\partial}{\partial w} [\ln f(w)] \Big|_{w=m}$$

Advantages: gives an approx distribution we can reason about.
Second derivatives are often tractable!

Limitations: bad if $p(w)$ is multimodal

bad if $p(w)$ has heavy tails
not symmetric about mode.

Derivation of Taylor approx to density $p(w)$ at $m = \arg \max_w p(w)$

$$(\log) p(w) = \log f(w) + \text{const}$$

$$= \ell(w) + \text{const}$$

define $\ell(w) = \log f(w)$

$$= \ell(m) + \underbrace{\ell'(m)(w-m)}_{-\frac{1}{2}\ell''(m)(w-m)^2} + \frac{1}{2} \ell''(m)(w-m)^2 + C$$

2nd order Taylor approx
to func ℓ at $w=m$

$$= -\frac{1}{2} \ell''(m)(w-m)^2 + C$$

m is mode, $\ell'(m) = 0$

$$= -\frac{1}{2} \beta (w-m)^2$$

$$\text{where } \beta = \frac{1}{\sigma^2} = -\ell''(m) \text{ (precision)}$$

$$\text{mean: } m = \arg \max_m \ell(w)$$

Laplace Approx in M-Dim

Idea: Approx: $q(w) = \text{MVNorm}(m, \Lambda^{-1})$
 ↓ mean ↓ precision matrix

- Set m to match the mode of f

$$m = \arg \max_{w \in \mathbb{R}^m} \ell(w)$$

- Set Λ to negative Hessian at mode

$$\Lambda = \nabla_w \nabla_w [-\ell(w)]|_{w=m}$$

Summary: A Laplace approximation 近似 -> 它先寻找一个众数 z_0 (估计), 然后计算那个众数位置上的Hessian matrix.

- 根据's Prior 贝叶斯原理, 我们知道如果后验分布是高斯的, 那么对于一个高斯的后验分布, Laplace approximation 是准确的!

Laplace approximation for Logistic regression

True posterior intractable, but known to Unstain via Bayes.

$$p(w | t_{1:N}) = \frac{p(t_{1:N}|w) \cdot p(w)}{p(t_{1:N})}$$

$$\begin{aligned} \Rightarrow \log p(w | t_{1:N}) &= \log p(t_{1:N} | w) + \log p(w) - \underbrace{\log p(t_{1:N})}_{C} \\ &= \log \text{MVNormPDF}(w | m_0, S_0) + \sum_{i=1}^N \text{Bern PDF}(t_i | \sigma(w^T \phi(x_i))) \end{aligned}$$

$$\log L(w)$$

Apply Laplace approx:

$$p(w | \{x_n, y_n\}_1^n) \approx N(w | m_{MAP}, S)$$

where $m_{MAP} = \arg \max_w L(w)$

$$S^{-1} = \nabla_w \nabla_w [-L(w)]|_{w=m_{MAP}}$$

- Using formula for the Hessian of MAP objective, we know:

$$\begin{aligned} S^{-1} &= S_0^{-1} + \Phi^T R(M_{MAP}) \Phi \\ &= \lambda I_m + \Phi^T R(M_{MAP}) \Phi \end{aligned}$$

Recall that $\mathbb{E}^T R \Sigma = \sum_{n=1}^N r(\text{MAP}, X_n) \phi(X_n) \phi(X_n)^T$

key point: Laplace parameters M, S^{-1} possible to compute!
predictive posteriors for Bayesian Logistic Regression

Ideal (intractable) posterior predictive:

$$p(t_* | t_{1:N}) = \int_{w \in \mathbb{R}^M} p(t_* | w) p(w | t_{1:N}) dw$$

Laplace approximation $\approx \int_{w \in \mathbb{R}^M} p(t_* | w) N(w | m_{\text{MAP}}, S) dw$

$$p(w | t_{1:N}) \approx N(m_{\text{MAP}}, S)$$

Still a tough integral over M dimensions
If M was 1 or 2, could use numerical strategies like trapezoid approx.

Option 1: Monte Carlo
Easy but need many samples L

$$p(t_* | t_{1:N}) = \mathbb{E}_{p(w | t_{1:N})} [p(t_* | w)]$$

Average of L samples from approx. posterior

$$\approx \frac{1}{L} \sum_{l=1}^L p(t_* | w^l)$$

with $w^l \sim N(m_{\text{MAP}}, S)$

$$\int [p(t_* | w, t_{1:N})] dw$$

$$= \int p(w | t_{1:N}) \cdot p(t_* | w) dw$$

$$= \mathbb{E}_{p(w | t_{1:N})} [p(t_* | w)]$$

Option 2: Probit approx.
Closed-form, hard to port to other models

Likelihood $p(t_* | w) = \begin{cases} \pi(w^T \phi_*) & \text{if } t_* = 1 \\ 1 - \pi(w^T \phi_*) & \text{if } t_* = 0 \end{cases}$

See Bishop Fig 4.9 $\approx \begin{cases} \text{NormCDF}(\sqrt{\frac{1}{S}} w^T \phi_*) & \text{if } t_* = 1 \\ 1 - \text{above} & \text{if } t_* = 0 \end{cases}$

$\pi(a) \approx \text{NormCDF}(\sqrt{\frac{1}{S}} a)$

Makes integral tractable when combined w/ Laplace!

$$(4.155): p(t_* = 1 | t_{1:N}) = \pi\left(\frac{m_{\text{MAP}}^T \phi_* + \frac{1}{\sqrt{1 + \frac{1}{S} \phi_*^T S \phi_*}}}{\sqrt{1 + \frac{1}{S} \phi_*^T S \phi_*}}\right)$$

