

$$\text{Model / Mappin: } y = f(\phi(x), \omega)$$

features : $\phi(x)$

ω : parameters (unknowns)

target : t

OBJECTIVE FUNCTION

$$J(\omega) = \frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2$$

$$\begin{aligned}
 & \underset{\omega}{\arg \min} J(\omega) \\
 &= \underset{\omega}{\arg \max} (-J(\omega)) \\
 &\stackrel{P}{=} \underset{\omega}{\arg \max} \exp(-J(\omega))
 \end{aligned}$$

$\exp(\cdot)$
is
monotonic

$$= \underset{\omega}{\arg \max} \exp\left(-\frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2\right)$$

\downarrow

$$= \underset{\omega}{\arg \max} \prod_{i=1}^N \underbrace{\exp\left(-\frac{1}{2} (t_i - y_i)^2\right)}_{\sim G(y_i, 1)}$$

$\exp(a+b) = \exp(a) \cdot \exp(b)$

$$= \underset{\omega}{\arg \max} \prod_{i=1}^N G(t_i; y_i, 1)$$

$L^0 = \text{observed data likelihood}$

Reminder

Gaussian

R.V.

$$x \sim G(\mu, \sigma^2)$$

↑ ↑
mean variance σ^2

PDF: $f_x(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}\right)$

OBSERVED
DATA

Likelihood

$$\mathcal{L}^o = \prod_{i=1}^N G(t_i; y_i, 1)$$

where $y_i = f(\omega, \phi(x_i))$

DATASET : $\{(x_i, t_i)\}_{i=1}^N$, ASSUME x_i are i.i.d.
independent and identically distributed

$$\mathcal{L}^o = P(x_1, x_2, x_3, \dots, x_N | \omega)$$

$$= P(x_1 | \omega) \cdot P(x_2 | \omega) \cdot \dots \cdot P(x_N | \omega)$$

\leftarrow
 x_i are
conditionally
independent

$$= \prod_{i=1}^N P(x_i | \omega)$$

$$= \prod_{i=1}^N G(t_i; y_i, 1)$$

Now adding a regularizer to obj-fct.

$$J(\omega) = \frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2 + \frac{\lambda}{2} \sum_j \omega_j^2$$

RIDGE
REG.

$$\arg \min_{\omega} J(\omega) = \arg \max_{\omega} (-J(\omega))$$

$$= \arg \max_{\omega} \exp(-J(\omega))$$

$$= \arg \max_{\omega} \exp\left(-\frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2 - \frac{\lambda}{2} \sum_j \omega_j^2\right)$$

$$= \arg \max_{\omega} \exp\left(-\frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2\right) \cdot \exp\left(-\frac{\lambda}{2} \sum_j \omega_j^2\right)$$

= \star

$$\hat{\omega} = \arg \max_{\omega} \left[\prod_{i=1}^N \underbrace{\exp \left(-\frac{1}{2} (t_i - y_i)^2 \right)}_{\sim G(t_i; y_i, 1)} \right] \left[\prod_j \underbrace{\exp \left(-\frac{\lambda}{2} \omega_j^2 \right)}_{\sim G(\omega_j; 0, 1/\lambda)} \right]$$

OBSERVED DATA
 Likelihood $\equiv \mathcal{L}^0$
 $\equiv P(t|\omega)$
Prior PROBABILITY
 i.e. probabilistic
 MODEL DESCRIBING
 THE UNKNOWN
 PARAMETERS ω
 $\equiv P(\omega)$

$$= \arg \max_{\omega} P(t|\omega) \cdot P(\omega)$$

$$= \arg \max_{\omega} P(\omega|t) \cdot P(t)$$

↙
 Bayes' Theorem

$$\propto \arg \max_{\omega} P(\omega|t)$$

Proportional

Reminder:

BAYES' THEOREM:

$P(\omega | t)$

posterior

$$P(\omega | t) = \frac{\underbrace{P(t | \omega)}_{\text{DATA likelihood}} \times \underbrace{P(\omega)}_{\text{prior prob}}}{\underbrace{P(t)}_{\text{Evidence prob.}}}$$

OBJ. FCT with LASSO REGULARIZER

$$J(\omega) = \frac{1}{2} \sum_{i=1}^N (t_i - y_i)^2 + \lambda \sum_j |\omega_j|$$

$$\arg \min_{\omega} J(\omega)$$

:

$$= \arg \max_{\omega} \left[\prod_{i=1}^N \underbrace{\exp \left(-\frac{1}{2} (t_i - y_i)^2 \right)}_{G(t_i; y_i, 1)} \right] \left[\prod_j \underbrace{\exp \left(-\lambda |\omega_j| \right)}_{L(\omega_j; 0, 1/\lambda)} \right]$$

reminder

Laplacian (μ, b)

$X \sim \text{Lap.}(\mu, b)$

$$f_X(x) = \frac{1}{2b} \cdot \exp\left(-\frac{|x-\mu|}{b}\right)$$

MAXIMUM LIKELIHOOD ESTIMATION (MLE)

- completely DATA - DRIVEN
- it requires a lot more data to make correct estimations for w

$$\arg \max_w P(t|w)$$

MAXIMUM A POSTERIORI (MAP)

- it requires a prior prob.
- prior will inject beliefs as to what values one most likely for w .
- if prior is selected correctly then w_{MAP} is better than w_{MLE} .
- if not, w_{MAP} will typically require a lot more data to compensate.

$$\arg \max_w P(t|w) \cdot P(w)$$

Example: flip a coin 3 times
and we observe $H_1 \cap H_2 \cap H_3 = E$

outcome : $E = H_1 \cap H_2 \cap H_3$

MLE (Frequentist approach) :

$$P(\text{Heads}) = \frac{\# \text{ Heads}}{\# \text{ flips}} = \frac{3}{3} = 1$$

MAP (Bayesian approach)

1st start with a set of prior beliefs. Beliefs: fair coin or 2-headed coin

$$P(\mu = \frac{1}{2} \mid \text{fair} \mid E)$$

$$P(\mu = 1 \mid E)$$

$$P(\mu = \frac{1}{2} \mid E) \stackrel{\mu = \frac{1}{2}}{\gtrsim} P(\mu = 1 \mid E)$$

$$P(\mu | E) = \frac{P(E|\mu) \times P(\mu)}{P(E)}$$

$$E = H_1 \cap H_2 \cap H_3$$

$$E = H_L$$

$$P(\mu | H) = \frac{P(H|\mu) \times P(\mu)}{P(H)}$$

$$\mu = \frac{1}{2} : = \frac{\frac{1}{2} \times \frac{1}{2}}{P(H|\mu = \frac{1}{2}) \times P(\mu = \frac{1}{2}) + P(H|\mu = 1) \times P(\mu = 1)}$$

Law of Total Probability

$$\begin{aligned} \text{Probability} &= \frac{\frac{1}{2} \times \frac{1}{2}}{\frac{1}{2} \times \frac{1}{2} + 1 \times \frac{1}{2}} \\ &= \frac{\frac{1}{4}}{\frac{1}{4} + \frac{1}{2}} = \frac{1}{3} \end{aligned}$$

$$\text{Similarly, } P(\mu = 1 | H) = \frac{2}{3}$$

Since $P(\mu = 1 | H) > P(\mu = \frac{1}{2} | H)$

then $\hat{\mu} = 1$ i.e. a 2-headed coin is most likely to have generated this outcome.

HEADS \rightarrow 1
TAILS \rightarrow 0

$$S = \{0, 1\}$$

μ = UNKNOWN parameter
prob. of flipping heads

x = input data (0 or 1)

$$P(x=1|\mu) = \mu$$

$$P(x=0|\mu) = 1-\mu$$

BERNOULLI R.V.

$$P(x|\mu) = \begin{cases} 1-\mu & , x=0 \\ \mu & , x=1 \\ 0 & , \text{otherwise} \end{cases}$$

$x \sim \text{Ber}(\mu)$

TRAINING
DATASET

$$: \{x_i\}_{i=1}^N, x_i \in \{0, 1\}$$

Starting with
MLE

STEP 1

OBSERVED DATA

Likelihood:

$$\mathcal{L}^o = P(x_1 \cap x_2 \cap \dots \cap x_N | \mu)$$

$$= P_1(x_1 | \mu) \cdot P_2(x_2 | \mu) \cdot \dots \cdot P_N(x_N | \mu)$$



x_i are
conditionally
independent
identically
distributed

$$= \prod_{i=1}^N P(x_i | \mu)$$

$$= \prod_{i=1}^N \mu^{x_i} \cdot (1-\mu)^{1-x_i}$$

assuming $P(\cdot)$
is a Bernoulli

Example : $\alpha = \{1, 0, 1\}$

$$\mathcal{L}^0 = \mu \cdot (1-\mu) \cdot \mu = \mu^2 \cdot (1-\mu)$$

$$\arg \max_{\mu} \mathcal{L}^0$$

$$= \arg \max_{\mu} \left[\prod_{i=1}^N \mu^{x_i} \cdot (1-\mu)^{1-x_i} \right]$$

STEP 2

Apply the \log to \mathcal{L}^0 .

$$\mathcal{L} = \ln(\mathcal{L}^0)$$

$$= \ln \left[\prod_{i=1}^n \mu^{x_i} \cdot (1-\mu)^{1-x_i} \right]$$

$$= \sum_{i=1}^n \ln(\mu^{x_i} \cdot (1-\mu)^{1-x_i})$$

$$= \sum_{i=1}^n (\ln(\mu^{x_i}) + \ln((1-\mu)^{1-x_i}))$$

$$= \sum_{i=1}^n (x_i \cdot \ln(\mu) + (1-x_i) \cdot \ln(1-\mu))$$

$$\mathcal{L} = \sum_{i=1}^N [x_i \cdot \ln(\mu) + (1-x_i) \cdot \ln(1-\mu)]$$

STEP 3

Solve for μ : $\frac{\partial \mathcal{L}}{\partial \mu} = 0$

$$\Leftrightarrow \sum_{i=1}^N \left[x_i \cdot \frac{1}{\mu} + (1-x_i) \cdot \left(\frac{-1}{1-\mu} \right) \right] = 0$$

$$\Leftrightarrow \sum_{i=1}^N (x_i(1-\mu) + \mu(x_i - 1)) = 0$$

$$\Leftrightarrow \sum_{i=1}^N (x_i - x_i\cancel{\mu} + \cancel{\mu x_i} - \mu) = 0$$

$$\Leftrightarrow \sum_{i=1}^N x_i - \underbrace{\sum_{i=1}^N \mu}_{= N \cdot \mu} = 0$$

$$\Leftrightarrow \boxed{\mu_{MLE} = \frac{1}{N} \sum_{i=1}^N x_i}$$

$\equiv \frac{\# \text{heads}}{\# \text{flips}}$
 \equiv relative freq.
 for outcome
 Heads

MAP APPROACH

1st assume a prior probability

E.g., BETA DISTRIBUTION as prior

$$\text{Beta}(x|\alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \cdot x^{\alpha-1} \cdot (1-x)^{\beta-1}$$

where $\Gamma(x) = (x-1)!$

and $\alpha, \beta > 0$

STEP 1 OBSERVED DATA LIKELIHOOD

$$\mathcal{L}^0 = P(x|\mu) \times P(\mu|\alpha, \beta)$$

$$= \left(\prod_{i=1}^N \mu^{x_i} (1-\mu)^{1-x_i} \right) \cdot \left(\frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \mu^{\alpha-1} (1-\mu)^{\beta-1} \right)$$

$$= \mu^{\sum_{i=1}^N x_i} \cdot (1-\mu)^{\sum_{i=1}^N (1-x_i)} \cdot \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \cdot \mu^{\alpha-1} (1-\mu)^{\beta-1}$$

$$= \frac{\Gamma(\alpha+\beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \cdot \mu^{\sum_{i=1}^N x_i + \alpha - 1} \cdot (1-\mu)^{N - \sum_{i=1}^N x_i + \beta - 1}$$

$$\mathcal{L} \propto \mu^{\sum_{i=1}^N x_i + \alpha - 1} \cdot (1-\mu)^{N - \sum_{i=1}^N x_i + \beta - 1}$$

STEP 2

$$\mathcal{L} = \ln(\mathcal{L}^0)$$

$$= \left(\sum_{i=1}^n x_i + \alpha - 1 \right) \cdot \ln(\mu) + \left(N - \sum_{i=1}^n x_i + \beta - 1 \right) \cdot \ln(1-\mu)$$

STEP 3

$$\frac{\partial \mathcal{L}}{\partial \mu} = 0 \Leftrightarrow \left(\sum_{i=1}^n x_i + \alpha - 1 \right) \cdot \frac{1}{\mu}$$

$$- \left(N - \sum_{i=1}^n x_i + \beta - 1 \right) \cdot \frac{1}{1-\mu} = 0$$

$$\Leftrightarrow \mu_{MAP} = \frac{\sum_{i=1}^N x_i + \alpha - 1}{N + \alpha + \beta - 2}$$

$N \equiv \# \text{ samples}$

In MAP

BETA: $P(\mu | \alpha, \beta) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha) \cdot \Gamma(\beta)} \cdot \mu^{\alpha-1} (1-\mu)^{\beta-1}$

Prior

Posterior $P(\mu | \alpha) \propto$

$$\mu^{\sum_{i=1}^n x_i + \alpha - 1} (1-\mu)^{N - \sum_{i=1}^n x_i + \beta - 1}$$

when the prior and the
posterior have the same

parametric form (minus constants)
they are said to have a
Conjugate Prior Relationship

This means that we can
use the posterior as
the new prior in an
online fashion:

Other conjugate prior relationships:

Likelihood

Prior

Posterior

① Gaussian \times Gaussian = Gaussian

② Bernoulli \times Beta \propto Beta

:

many more.

For a new batch of data
iteration (t):

$$\alpha^{(t+1)} \leftarrow \alpha^{(t)} + \sum_{i=1}^n x_i$$

$$\beta^{(t+1)} \leftarrow \beta^{(t)} + n - \sum_{i=1}^n x_i$$

Review:

Multivariate Gaussian

Bivariate Gaussian

↳ covariance