

here we introduce another method: **Ridge regression**.

In the bayesian world, you would need to infer β parameters values, which you would then not go for the ML method (frequentist approach). Instead you would ask what is the most probable, likely β , that is most probable for seeing your data?

$$Z = \left\{ (x_i, y_i) \quad 1 \leq i \leq n \right\}, \text{ Full set of observations}$$

$$= \left\{ z_1, \dots, z_n \quad 1 \leq i \leq n \right\}$$

Bayes Formula: $P(\beta|Z) = \frac{\text{likelihood fn } \underbrace{P(Z|\beta)}_{\text{evidence}} \cdot \underbrace{P(\beta)}_{\text{prior}}}{P(Z)}$

Then you might go for: **maximum a posteriori estimate (MAP)**

$$\hat{\beta}_{\text{MAP}} \in \arg \max_{\beta} P(\beta|Z) \quad \text{posterior probability } P(\beta|Z).$$

$$= \arg \min_{\beta} \left\{ \underbrace{-\log P(Z|\beta)}_{\text{likelihood}} - \log P(\beta) \right\} \quad \text{minimizing ~~likelihood~~. chosen model}$$

* maximization of a likelihood fn under conditional gaussian noise distribution for a linear model is the same to minimize a sum-of-squares error fn.

REGULARIZATION

Before you look at the data, know that we assumed β to be distributed by nature according to gaussian law.

How to choose $P(\beta)$? prior

β is distributed by nature,

$$P(\beta) = \prod_{i=1}^d \frac{1}{\sqrt{2\pi\sigma_{\beta}^2}} \exp\left(-\frac{\beta_i^2}{2\sigma_{\beta}^2}\right)$$

Take the log of $P(\beta)$.

$$-\log P(\beta) = \sum_{j=1}^d \frac{\beta_j^2}{2\sigma_{\beta}^2} + \text{constant}$$

$d = \dim$
 $n = \text{examples}$

$$= \arg \min_{\beta} \left\{ \underbrace{\sum_{i=1}^n (y_i - x_i^T \beta)^2}_{\text{prediction error term}} + \underbrace{\lambda \sum_{j=1}^d \beta_j^2}_{\text{penalty regularizer, scalar to penalize max-likelihood fn}} \right\}$$

// cost fn for ridge regression

used to control model complexity, it's an isotropic type of penalty term for choosing the number of β 's you want to restrict.

cost fn depends on β .

$$RSS(\beta) = \underbrace{(Y - \hat{X}\beta)^T (Y - X\beta)}_{\text{data term}} + \underbrace{\lambda \|\beta\|^2}_{\text{penalty term}}$$

where $\lambda = \text{range parameter}$

$\lambda = \frac{1}{2\sigma_{\beta}^2}$, validated by data and controls solution to fitting parameters of the model.