Machine Learning and Computational Statistics Bayesian Method and Regression

Intro: This document consists of concepts and exercises related to Bayesian Method and Regression. The core idea lies in the believe that distribution of parameters changes and we use Bayes' theorem to update the probability for a hypothesis (e.g., parameter distribution) as more evidence or information becomes available.

Learning objectives of Bayesian Method and Regression include:

- 1. Basic Bayesian setup (i.e., Posterior distribution = Prior distribution * Likelihood)
- 2. The prior predictive distribution:

$$p(y|x) = \int p(y|x;\theta)p(\theta)d\theta$$

3. The posterior predictive distribution:

$$p(y|x, D) = \int p(y|x; \theta)p(\theta|D)d\theta$$

- 4. The difference between the posterior predictive distribution function (Bayesian approach) and the MAP (Posterior Mean Estimate) (Frequentist approach)
 - (a) Posterior predictive distribution does not depend on the unknown parameter because it has been integrated out.
 - (b) Maximum a posteriori probability MAP: we choose $\hat{\theta}$ and predict $p(y|x, \hat{\theta}(D))$. We then take derivative of the log-likelihood and set to 0.

$$\begin{split} P(\theta|D) &= \frac{P(D|\theta) * P(\theta)}{P(D)} \propto P(D|\theta) * P(\theta) \\ \hat{\theta}_{MAP} &= \underset{\theta}{\operatorname{argmax}} P(D|\theta) * P(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} log P(D|\theta) * P(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} log \prod_{i=1}^{n} P(y_i|x_i, w) * P(\theta) \\ &= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^{n} log P(y_i|x_i, w) * P(\theta) \end{split}$$

5. Conjugate prior:

- (a) Let π be a family of prior distributions on Θ
- (b) Let P be parametric family of distributions with parameter space Θ . P could be poisson, normal, or bernoulli
- (c) π is conjugate to P if for any prior in π , the posterior is always in π

6. Bayesian point estimates

- (a) posterior mean $\hat{\theta} = E(\theta|D)$
- (b) maximum a posteriori (MAP) estimator the mode of the posterior distribution

7.

Bayesian Decision Theory

- Ingredients:
 - Parameter space Θ.
 - **Prior**: Distribution $p(\theta)$ on Θ .
 - Action space A.
 - Loss function: $\ell: \mathcal{A} \times \Theta \to \mathbf{R}$.
- The **posterior risk** of an action $a \in A$ is

$$\begin{aligned} r(a) &:= & \mathbb{E}\left[\ell(\theta, a) \mid \mathcal{D}\right] \\ &= & \int \ell(\theta, a) p(\theta \mid \mathcal{D}) \, d\theta. \end{aligned}$$

- It's the expected loss under the posterior.
- A Bayes action a^* is an action that minimizes posterior risk:

$$r(a^*) = \min_{a \in \mathcal{A}} r(a)$$

8. Relationship between Gaussian regression and Ridge regression

Closed Form for Posterior

Model:

$$w \sim \mathcal{N}(0, \Sigma_0)$$

 $y_i \mid x, w \text{ i.i.d. } \mathcal{N}(w^T x_i, \sigma^2)$

- Design matrix X Response column vector y
- Posterior distribution is a Gaussian distribution:

$$\begin{array}{rcl} w \mid \mathcal{D} & \sim & \mathcal{N}(\mu_P, \Sigma_P) \\ \mu_P & = & \left(X^T X + \sigma^2 \Sigma_0^{-1} \right)^{-1} X^T y \\ \Sigma_P & = & \left(\sigma^{-2} X^T X + \Sigma_0^{-1} \right)^{-1} \end{array}$$

• Posterior Variance Σ_P gives us a natural uncertainty measure.

Closed Form for Posterior

Posterior distribution is a Gaussian distribution:

$$\begin{array}{rcl} \textbf{w} \mid \mathfrak{D} & \sim & \mathfrak{N}(\mu_P, \Sigma_P) \\ \mu_{\textbf{P}} & = & \left(\textbf{X}^T \textbf{X} + \sigma^2 \Sigma_0^{-1} \right)^{-1} \textbf{X}^T \textbf{y} \\ \Sigma_{\textbf{P}} & = & \left(\sigma^{-2} \textbf{X}^T \textbf{X} + \Sigma_0^{-1} \right)^{-1} \end{array}$$

• If we want point estimates of w, MAP estimator and the posterior mean are given by

$$\hat{w} = \mu_P = (X^T X + \sigma^2 \Sigma_0^{-1})^{-1} X^T y$$

• For the prior variance $\Sigma_0 = \frac{\sigma^2}{\lambda} I$, we get

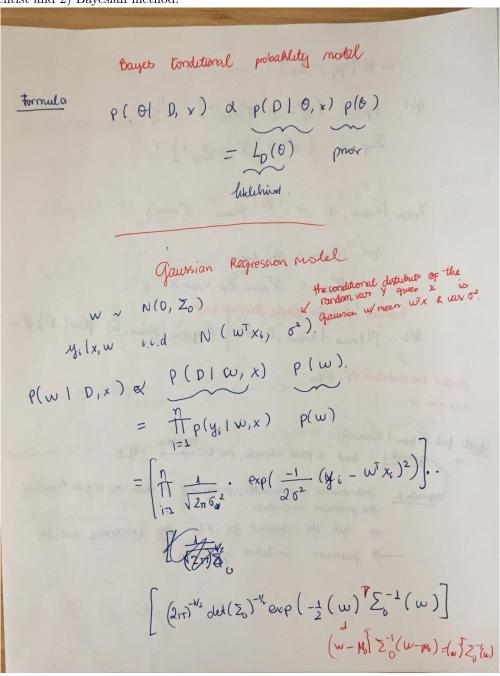
$$\hat{\mathbf{w}} = \mathbf{\mu}_P = \left(X^T X + \lambda I\right)^{-1} X^T \mathbf{y},$$

which is of course the ridge regression solution.

- 9. Concept check questions Bayesian Method and Regression
- 10. Bayesian Methods

1 Gaussian Regression

This Gaussian regression model is one example of how we can find y from two approaches of 1) Frequentist and 2) Bayesian method.



```
~ N (Mp, Zp)
       for Mp = (XTX + 82 \sum_{0}^{-1}) -1 XTY
             \Sigma_{\rho} = \left( \vec{\sigma}^2 \times^{\top} \times + \Sigma_{\rho}^{-1} \right)^{-1}
       Ynew Ixnew, O ~ N (Mnew, Orner)
              Not More = Mot xnew
                  Frew Ep xnow + 52
          Posterior predictive prindro distribute finch
        b/c: p(ynew (xnew, D) = ) p(ynew (xnew, w) p(w (D) du
  predict the dustribut OF Your
  aguer X
Goal: frid P(Ynow 1 xnew, D).
      approach! find a point estimale for W >>> MLE
      approach 2: unknown w is a variable, producing a distr on WEIR & called
               the posterior distributo
              => get the dustribut for YIX by integraty out W
           -> posterior predictive junction
```

Frequential approach
$$\rightarrow$$
 MLE to ind ω^*

Y!x \sim N ($\omega^T x$, σ^2). Y

Pw(y; |x|) = $\frac{1}{\sigma^{1/2}\pi}$ $\exp\left(\frac{1}{2\sigma^2}(y_i - \omega^T x_i)^2\right) \rightarrow \text{gath p(Ybrand)}$

L = PN(D)= $\prod_{i=1}^{n} \frac{1}{\sigma^{1/2}\pi} \exp\left(\frac{1}{2\sigma^2}(y_i - \omega^T x_i)^2\right)$.

log L = $\sum_{i=1}^{n} [\log_{\sigma^{1/2}\pi} + (\frac{1}{2\sigma^2}(y_i - \omega^T x_i)^2)]$.

Need to profess where $\omega^* = \alpha \log_{\sigma^{1/2}\pi} \log_{\sigma^{1/2}\pi} (y_i - \omega^T x_i)^2$

Need to profess where $\omega^* = \omega^* \log_{\sigma^{1/2}\pi} \log_{\sigma^{1/2}\pi} (y_i - \omega^T x_i)^2$
 $\frac{\partial P_{V}(D)}{\partial \omega} = -\frac{1}{2\sigma^2} \cdot \sum_{i=1}^{n} (y_i - \omega^T x_i) (-x_i)$
 $\frac{\partial P_{V}(D)}{\partial \omega} = -\frac{1}{2\sigma^2} \cdot \sum_{i=1}^{n} (y_i - \omega^T x_i) (-x_i)$
 $\frac{\partial P_{V}(D)}{\partial \omega} = \frac{1}{\sigma^2} \sum_{i=1}^{n} (y_i - \omega^T x_i) (-x_i)$
 $\frac{\partial P_{V}(D)}{\partial \omega} = \frac{1}{\sigma^2} \sum_{i=1}^{n} (y_i - \omega^T x_i) (-x_i)$

A: $\sum_{i=1}^{n} (y_i - w^i x_i)^2 = (y - xw)^T (y - xw)$ for $X = \begin{bmatrix} -x_{1-} \\ -x_{n-} \end{bmatrix}$ $= Y^{T}Y - 2w^{T}X^{T}Y + w^{T}X^{T}Xw.$ $\frac{\partial w cx}{\partial x} = c^{T} \frac{\partial x (Ax)}{\partial x} = (A^{T} + A)x.$ $\frac{\partial w}{\partial w} = -2x^{T}Y + \left[(x^{T}X)^{T} + (x^{T}X)\right]w$ $-2x^{T}y + 2x^{T}xw$ $= \frac{1}{X^{T}XW} = \frac{X^{T}Y}{X^{T}Y}$ P Note $X^{T}X$ may not be invertible!!!

B. Boyesian Approach

W is treated as a val.

$$P(w) (\text{ prior distribute of } w) \sim N(0, \Sigma).$$

$$f(x_i) = w^T x_i$$

$$E_i \sim N(0, \delta^1)$$

$$y_i = f(x_i) + E_i$$

The matrix form
$$= w \sim N(0, \Sigma).$$

$$f(X) = Xw$$

$$E \sim N(0, \delta^2 I)$$

$$Y = f(X) + E$$

$$= 1 \text{ Y | X,w} \sim N(Xw, \delta^2 I) = \text{All | X,w} \text{ And | } \text{ And$$

$$\frac{1}{2} P(y, | x_{1}, \omega) = \frac{1}{22\pi^{2}} \cdot exp(\frac{1}{20^{2}} (y_{1} - \omega^{T} x_{1})^{2}).$$
We with
$$P(D | x_{1}, \omega) = \prod_{i=1}^{T} P(y_{i}, x_{1}, \omega).$$

$$\frac{1}{2} exp(\frac{1}{20^{2}} \cdot \sum_{i=1}^{T} (y_{i} - \omega^{T} x_{i})^{2}).$$

$$= exp(\frac{1}{20^{2}} \cdot \sum_{i=1}^{T} (y_{i} - \omega^{T} x_{i})^{2}).$$

$$= exp(\frac{1}{20^{2}} \cdot (y - x\omega)^{T} (y - x\omega).$$

$$= exp(\frac{1}{20^{2}} \cdot (y - x\omega)^{T} (y - x\omega).$$

$$= w^{T} (\frac{1}{6^{2}} x^{T} x + \sum_{i=1}^{T} w - 2 \frac{1}{6^{2}} (y^{T} x w) + \frac{1}{6^{2}} y^{T} y.$$

$$= w^{T} (\frac{1}{6^{2}} x^{T} x + \sum_{i=1}^{T} w - 2 \frac{1}{6^{2}} (y^{T} x w) + \frac{1}{6^{2}} y^{T} y.$$

$$= w^{T} (w - 2 + w) + w^{T} x^{T} y.$$

$$= w^{T} (w - 2 + w) + w^{T} x^{T} y.$$

$$= w^{T} (w - 2 + w) + w^{T} x^{T} y.$$

$$= w^{T} (w - 2 + w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$

$$= w^{T} (w - w) + w^{T} x^{T} y.$$