

CPSC 540: Machine Learning

Empirical Bayes

Mark Schmidt

University of British Columbia

Winter 2020

Last Time: Bayesian Statistics

- For most of the course, we considered **MAP estimation**:

$$\hat{w} = \operatorname{argmax}_w p(w | X, y) \quad (\text{train})$$

$$\hat{y}^i = \operatorname{argmax}_{\hat{y}} p(\hat{y} | \hat{x}^i, \hat{w}) \quad (\text{test}).$$

- But w was random: I have **no justification** to only base decision on \hat{w} .
 - Ignores other reasonable values of w that could make opposite decision.
- Last time we introduced **Bayesian** approach:
 - Treat w as a **random variable**, and **define probability over what we want** given data:

$$\begin{aligned} \hat{y}^i &= \operatorname{argmax}_{\hat{y}} p(\hat{y} | \hat{x}^i, X, y) \\ &= \operatorname{argmax}_{\hat{y}} \int_w p(\hat{y} | \hat{x}^i, w) p(w | X, y) dw. \end{aligned}$$

- Directly follows from rules of probability, and no separate training/testing.

Bayesian Linear Regression

- We know that L2-regularized linear regression,

$$\operatorname{argmin}_w \frac{1}{2\sigma^2} \|Xw - y\|^2 + \frac{\lambda}{2} \|w\|^2,$$

corresponds to MAP estimation in the model

$$y^i \sim \mathcal{N}(w^T x^i, \sigma^2), \quad w_j \sim \mathcal{N}(0, \lambda^{-1}).$$

- By some tedious Gaussian identities, the posterior has the form

$$w \mid X, y \sim \mathcal{N} \left(\frac{1}{\sigma^2} \left(\frac{1}{\sigma^2} X^T X + \lambda I \right)^{-1} X^T y, \left(\frac{1}{\sigma^2} X^T X + \lambda I \right)^{-1} \right).$$

- Notice that mean of posterior is the MAP estimate (not true in general).
- Bayesian perspective gives us variability in w and optimal predictions given prior.
- But it also gives different ways to choose λ and choose basis.

Learning the Prior from Data?

- Can we use the training data to set the hyper-parameters?
- In theory: No!
 - It would not be a “prior”.
 - It’s no longer the right thing to do.
- In practice: Yes!
 - Approach 1: split into training/validation set or use cross-validation as before.
 - Approach 2: optimize the **marginal likelihood** (“evidence”):

$$p(y | X, \lambda) = \int_w p(y | X, w)p(w | \lambda)dw.$$

- Also called **type II maximum likelihood** or **evidence maximization** or **empirical Bayes**.

Digression: Marginal Likelihood in Gaussian-Gaussian Model

- Suppose we have a **Gaussian likelihood** and **Gaussian prior**,

$$y^i \sim \mathcal{N}(w^T x^i, \sigma^2), \quad w_j \sim \mathcal{N}(0, \lambda^{-1}).$$

- The joint probability of y^i and w_j is given by

$$p(y, w | X, \lambda) \propto \exp\left(-\frac{1}{2\sigma^2}\|Xw - y\|^2 - \frac{\lambda}{2}\|w\|^2\right).$$

- The **marginal likelihood integrates** the joint over the nuisance parameter w ,

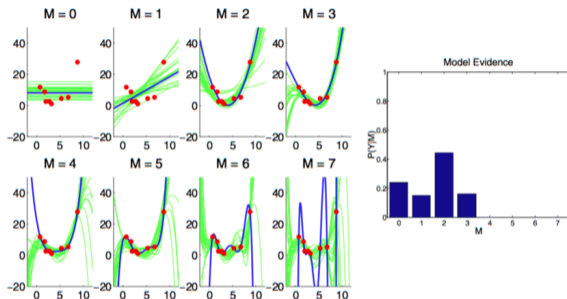
$$p(y | X, \lambda) = \int_w p(y, w | X, \lambda) dw.$$

- Solving the Gaussian integral gives a **marginal likelihood** of

$$p(y | X, \lambda) \propto |C|^{-1/2} \exp\left(-\frac{y^T C^{-1} y}{2}\right), \quad C = \sigma^2 I + \frac{1}{\lambda} X X^T.$$

Type II Maximum Likelihood for Basis Parameter

- Consider **polynomial basis**, and treat degree M as a hyper-parameter:



http://www.cs.ubc.ca/~arnaud/stat535/slides5_revised.pdf

- Marginal likelihood (evidence) is highest for $M = 2$.
 - "Bayesian Occam's Razor": prefers simpler models that fit data well.
 - $p(y | X, \lambda)$ is small for $M = 7$, since 7-degree polynomials can fit many datasets.
 - It's actually **non-monotonic** in M : it prefers $M = 0$ and $M = 2$ over $M = 1$.
 - Model selection criteria like BIC are approximations to marginal likelihood as $n \rightarrow \infty$.

Type II Maximum Likelihood for Polynomial Basis

- Why is the marginal likelihood **high for degree 2 but not degree 7**?
 - Marginal likelihood for degree 2:

$$p(y | X, \lambda) = \int_{w_0} \int_{w_1} \int_{w_2} p(y | X, w) p(w | \lambda) dw$$

- Marginal likelihood for degree 7:

$$p(y | X, \lambda) = \int_{w_0} \int_{w_1} \int_{w_2} \int_{w_3} \int_{w_4} \int_{w_5} \int_{w_6} \int_{w_7} p(y | X, w) p(w | \lambda) dw.$$

- Higher-degree integrates over high-dimensional volume:
 - A non-trivial **proportion** of degree 2 functions fit the data really well.
 - There are many degree 7 functions that fit the data even better, but they are a **much smaller proportion** of all degree 7 functions.
- Warning: this doesn't always work, sometimes becomes degenerate.
 - May **need a prior on the hyper-parameters**.

Bayes Factors for Bayesian Hypothesis Testing

- Suppose we want to **compare hypotheses**:
 - E.g., “this data is best fit with linear model” vs. a degree-2 polynomial.
- **Bayes factor** is ratio of marginal likelihoods,

$$\frac{p(y | X, \text{degree } 2)}{p(y | X, \text{degree } 1)}$$

- If very large then data is much more consistent with degree 2.
 - A common variation also puts **prior on degree**.
- A more **direct method of hypothesis testing**:
 - No need for null hypothesis, “power” of test, p-values, and so on.
 - As usual can only tell you which model is likely, not whether any are correct.

- American Statistical Association:
 - “Statement on Statistical Significance and P-Values” .
 - <http://amstat.tandfonline.com/doi/pdf/10.1080/00031305.2016.1154108>
- “Hack Your Way To Scientific Glory” :
 - <https://fivethirtyeight.com/features/science-isnt-broken>
- “Replicability crisis” in social psychology and many other fields:
 - https://en.wikipedia.org/wiki/Replication_crisis
 - <http://www.nature.com/news/big-names-in-statistics-want-to-shake-up-much-maligned-p-value-1.22375>
- “T-Tests Aren't Monotonic” : <https://www.naftaliharris.com/blog/t-test-non-monotonic>
- Bayes factors don't solve problems with p-values and multiple testing.
 - But they give an alternative view, are more intuitive, and make assumptions clear.
- Some notes on various issues associated with Bayes factors:
 - <http://www.aarondefazio.com/aderazio-bayesfactor-guide.pdf>

Learning Principles

- Maximum likelihood:

$$\hat{w} \in \operatorname{argmax}_w p(y | X, w) \qquad \hat{y} \in \operatorname{argmax}_{\tilde{y}} p(\tilde{y} | \tilde{x}, \hat{w}).$$

- MAP:

$$\hat{w} \in \operatorname{argmax}_w p(w | X, y, \lambda) \qquad \hat{y} \in \operatorname{argmax}_{\tilde{y}} p(\tilde{y} | \tilde{x}, \hat{w}).$$

- Optimizing λ in this setting **does not work**: sets $\lambda = 0$.
- Bayesian (no “learning”):

$$\hat{y} \in \operatorname{argmax}_{\tilde{y}} \int_w p(\tilde{y} | \tilde{x}, w) p(w | X, y, \lambda) dw.$$

- Type II maximum likelihood (“learn hyper-parameters”):

$$\hat{\lambda} \in \operatorname{argmax}_{\lambda} p(y | X, \lambda) \qquad \hat{y} \in \operatorname{argmax}_{\tilde{y}} \int_w p(\tilde{y} | \tilde{x}, w) p(w | X, y, \hat{\lambda}) dw.$$

Type II Maximum Likelihood for Regularization Parameter

- Type II maximum likelihood maximizes probability of data given hyper-parameters,

$$\hat{\lambda} \in \underset{\lambda}{\operatorname{argmax}} p(y | X, \lambda), \quad \text{where} \quad p(y | X, \lambda) = \int_w p(y | X, w)p(w | \lambda)dw,$$

and the integral has closed-form solution if everything is Gaussian.

- You can run gradient descent to choose λ .
- We are using the data to optimize the prior (empirical Bayes).
- Even if we have a complicated model, much less likely to overfit than MLE:
 - Complicated models need to integrate over many more alternative hypotheses.

Type II Maximum Likelihood for Individual Regularization Parameter

- Consider having a hyper-parameter λ_j for each w_j ,

$$y^i \sim \mathcal{N}(w^T x^i, \sigma^2 I), \quad w_j \sim \mathcal{N}(0, \lambda_j^{-1}).$$

- Too expensive for cross-validation, but type II MLE works.
 - You can do **gradient descent to optimize the λ_j** .
- Weird fact: this yields **sparse** solutions.
 - “**Automatic relevance determination**” (ARD)
 - Can send $\lambda_j \rightarrow \infty$, concentrating posterior for w_j at exactly 0.
 - It tries to “remove some of the integrals”.
 - This is L2-regularization, but **empirical Bayes naturally encourages sparsity**.
- Non-convex and theory not well understood:
 - Tends to yield much sparser solutions than L1-regularization.

Type II Maximum Likelihood for Other Hyper-Parameters

- Consider also having a hyper-parameter σ_i for each i ,

$$y^i \sim \mathcal{N}(w^T x^i, \sigma_i^2), \quad w_j \sim \mathcal{N}(0, \lambda_j^{-1}).$$

- You can also use type II MLE to optimize these values.
- The “automatic relevance determination” selects training examples ($\sigma_i \rightarrow \infty$).
 - This is like the support vectors in SVMs, but tends to be much more sparse.
- Type II MLE can also be used to learn kernel parameters like RBF variance.
 - Do gradient descent on the σ values in the Gaussian kernel.
- It may also do something sensible if you use it to choose number of clusters k .
 - Or number of states in hidden Markov model, number of latent factors in PCA, etc.
- Bonus slides: Bayesian feature selection gives probability that w_j is non-zero.
 - Posterior is much more informative than standard sparse MAP methods.

Summary

- **Marginal likelihood** is probability seeing data given hyper-parameters.
- **Empirical Bayes** optimizes marginal likelihood to set hyper-parameters:
 - Allows tuning a large number of hyper-parameters.
 - Bayesian Occam's razor: naturally encourages sparsity and simplicity.
- Next time: which priors yield closed-form solutions?

Gradient on Validation/Cross-Validation Error

- It's also possible to do **gradient descent on λ to optimize validation/cross-validation error** of model fit on the training data.
- For L2-regularized least squares, define $w(\lambda) = (X^T X + \lambda I)^{-1} X^T y$.
- You can use chain rule to get **derivative of validation error E_{valid} with respect to λ** :

$$\frac{d}{d\lambda} E_{\text{valid}}(w(\lambda)) = E'_{\text{valid}}(w(\lambda)) w'(\lambda).$$

- For more complicated models, you can use **total derivative** to get gradient with respect to λ in terms of gradient/Hessian with respect to w .
- However, this is often more sensitive to over-fitting than empirical Bayes approach.

Bayesian Feature Selection

- Classic feature selection methods don't work when $d \gg n$:
 - AIC, BIC, Mallows', adjusted- R^2 , and L1-regularization return very different results.
- Here maybe all we can hope for is **posterior probability of $w_j = 0$** .
 - Consider all models, and weight by posterior the ones where $w_j = 0$.
- If we fix λ and use L1-regularization, posterior is **not sparse**.
 - Probability that a variable is exactly 0 is zero.
 - L1-regularization only leads to sparse MAP, not sparse posterior.

Bayesian Feature Selection

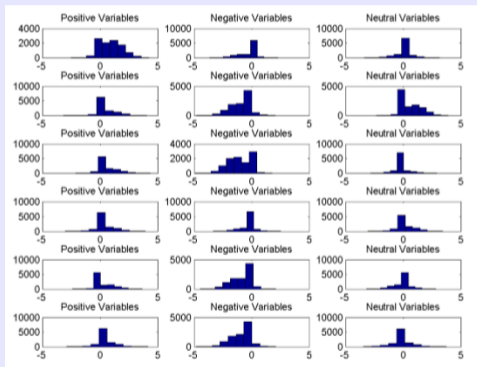
- Type II MLE gives sparsity because posterior variance goes to zero.
 - But this **doesn't give probability** of being 0.
- We can encourage sparsity in Bayesian models using a **spike and slab** prior:



- Mixture of Dirac delta function at 0 and another prior with non-zero variance.
- Places non-zero posterior weight at exactly 0.
- Posterior is still non-sparse, but answers the question:
 - “What is the probability that variable is non-zero”?

Bayesian Feature Selection

- Monte Carlo samples of w_j for 18 features when classifying '2' vs. '3':
 - Requires “trans-dimensional” MCMC since dimension of w is changing.



- “Positive” variables had $w_j > 0$ when fit with L1-regularization.
- “Negative” variables had $w_j < 0$ when fit with L1-regularization.
- “Neutral” variables had $w_j = 0$ when fit with L1-regularization.