STA 314: Statistical Methods for Machine Learning I Lecture 11 - Bayesian Linear Regression, Probabilistic PCA

Chris J. Maddison

University of Toronto

- Final exam does not include this nor next week's lecture.
- Continuing in our theme of probabilistic models for continuous variables.
 - Probabilistic interpretation of linear regression
 - Probabilistic interpretation of PCA
- (Optional) Bayesian model selection.

Completing the Square for Gaussians

- First, we're going to review a very powerful technique that will let us figure out the distribution of Gaussian random variables.
- It's a multivariate generalization of completing the square.
- The density of $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ satifies:

$$\log p(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^{\mathsf{T}} \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) + \text{const}$$
$$= -\frac{1}{2} \mathbf{x}^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \mathbf{x} + \mathbf{x}^{\mathsf{T}} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu} + \text{const}$$

• Thus, if we know **w** is Gaussian with *unknown* mean and covariance, and we also know that

$$\log p(\mathbf{w}) = -\frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathbf{A}\mathbf{w} + \mathbf{w}^{\mathsf{T}}\mathbf{b} + \text{const}$$

for A positive definite, then we know that

$$\mathbf{w} \sim \mathcal{N}(\mathbf{A}^{-1}\mathbf{b}, \mathbf{A}^{-1})$$

• We're going to be Bayesian about the parameters of the model.

- This is in contrast with naïve Bayes and GDA: in those cases, we used Bayes' rule to infer the class, but used point estimates of the parameters.
- By inferring a posterior distribution over the *parameters*, the model can know what it doesn't know.
- How can uncertainty in the predictions help us?
 - Smooth out the predictions by averaging over lots of plausible explanations (just like ensembles!)
 - Assign confidences to predictions
 - Make more robust decisions

Recap: Linear Regression

- Given a training set of inputs and targets $\{(\mathbf{x}^{(i)}, t^{(i)})\}_{i=1}^{N}$
- Linear model:

$$y = \mathbf{w}^{\top} \psi(\mathbf{x})$$

 \bullet Vectorized, we have the design matrix ${\boldsymbol X}$ in input space and

$$\Psi = \begin{bmatrix} - & \psi(\mathbf{x}^{(1)}) & - \\ - & \psi(\mathbf{x}^{(2)}) & - \\ \vdots & \\ - & \psi(\mathbf{x}^{(N)}) & - \end{bmatrix}$$

and predictions

Recap: Linear Regression

Squared error loss:

$$L(\mathbf{y},\mathbf{t}) = \frac{1}{2} \|\mathbf{y} - \mathbf{t}\|^2$$

• L₂ regularization:

$$\phi(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|^2$$

• Solution 1: solve analytically by setting the gradient to 0

$$\mathbf{w} = (\mathbf{\Psi}^{\top} \mathbf{\Psi} + \lambda \mathbf{I})^{-1} \mathbf{\Psi}^{\top} \mathbf{t}$$

• Solution 2: solve approximately using gradient descent

$$\mathbf{w} \leftarrow (1 - \alpha \lambda) \mathbf{w} - \alpha \mathbf{\Psi}^{\mathsf{T}} (\mathbf{y} - \mathbf{t})$$

Linear Regression as Maximum Likelihood

• We can give linear regression a probabilistic interpretation by assuming a Gaussian noise model:

$$t \mid \mathbf{x} \sim \mathcal{N}(\mathbf{w}^{\top} \boldsymbol{\psi}(\mathbf{x}), \sigma^2)$$

• Linear regression is just maximum likelihood under this model:

$$\frac{1}{N} \sum_{i=1}^{N} \log p(t^{(i)} | \mathbf{x}^{(i)}; \mathbf{w}, b) = \frac{1}{N} \sum_{i=1}^{N} \log \mathcal{N}(t^{(i)}; \mathbf{w}^{\top} \boldsymbol{\psi}(\mathbf{x}^{(i)}), \sigma^{2})$$
$$= \frac{1}{N} \sum_{i=1}^{N} \log \left[\frac{1}{\sqrt{2\pi\sigma}} \exp\left(-\frac{(t^{(i)} - \mathbf{w}^{\top} \boldsymbol{\psi}(\mathbf{x}^{(i)}))^{2}}{2\sigma^{2}}\right) \right]$$
$$= \operatorname{const} - \frac{1}{2N\sigma^{2}} \sum_{i=1}^{N} (t^{(i)} - \mathbf{w}^{\top} \boldsymbol{\psi}(\mathbf{x}^{(i)}))^{2}$$

Regularized Linear Regression as MAP Estimation

- We can view an L_2 regularizer as MAP inference with a Gaussian prior.
- Recall MAP inference:

$$\arg\max_{\mathbf{w}} \log p(\mathbf{w} \mid \mathcal{D}) = \arg\max_{\mathbf{w}} [\log p(\mathbf{w}) + \log p(\mathcal{D} \mid \mathbf{w})]$$

• We just derived the likelihood term $\log p(\mathcal{D} | \mathbf{w})$:

$$\log p(\mathcal{D} \mid \mathbf{w}) = \text{const} - \frac{1}{2N\sigma^2} \sum_{i=1}^{N} (t^{(i)} - \mathbf{w}^{\top} \psi(\mathbf{x}^{(i)}))^2$$

• Assume a Gaussian prior, $\mathbf{w} \sim \mathcal{N}(\mathbf{m}, \mathbf{S})$:

$$\log p(\mathbf{w}) = \log \mathcal{N}(\mathbf{w}; \mathbf{m}, \mathbf{S})$$

=
$$\log \left[\frac{1}{(2\pi)^{D/2} |\mathbf{S}|^{1/2}} \exp \left(-\frac{1}{2} (\mathbf{w} - \mathbf{m})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{w} - \mathbf{m}) \right) \right]$$

=
$$-\frac{1}{2} (\mathbf{w} - \mathbf{m})^{\mathsf{T}} \mathbf{S}^{-1} (\mathbf{w} - \mathbf{m}) + \text{const}$$

• Commonly, $\mathbf{m} = \mathbf{0}$ and $\mathbf{S} = \eta \mathbf{I}$, so

$$\log p(\mathbf{w}) = -\frac{1}{2\eta} ||\mathbf{w}||^2 + \text{const.}$$

This is just L_2 regularization!

Intro ML (UofT)

- Recall: full Bayesian inference makes predictions by averaging over all likely explanations under the posterior distribution.
- Compute posterior using Bayes' Rule:

$$p(\mathbf{w} \mid \mathcal{D}) \propto p(\mathbf{w})p(\mathcal{D} \mid \mathbf{w})$$

• Make predictions using the posterior predictive distribution:

$$p(t \mid \mathbf{x}, \mathcal{D}) = \int p(\mathbf{w} \mid \mathcal{D}) p(t \mid \mathbf{x}, \mathbf{w}) \, \mathrm{d}\mathbf{w}$$

• Doing this lets us quantify our uncertainty.

- Prior distribution: w ~ $\mathcal{N}(\mathbf{0}, \mathbf{S})$
- Likelihood: $t \mid \mathbf{x}, \mathbf{w} \sim \mathcal{N}(\mathbf{w}^{\top} \boldsymbol{\psi}(\mathbf{x}), \sigma^2)$
- \bullet Assuming fixed/known ${\bf S}$ and σ^2 is a big assumption. More on this later.

- Bayesian linear regression considers various plausible explanations for how the data were generated.
- It makes predictions using all possible regression weights, weighted by their posterior probability.
- Here are samples from the prior $p(\mathbf{w})$ and posteriors $p(\mathbf{w} | D)$



• Deriving the posterior distribution:

$$\log p(\mathbf{w} \mid \mathcal{D}) = \log p(\mathbf{w}) + \log p(\mathcal{D} \mid \mathbf{w}) + \text{const}$$

$$= -\frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathbf{S}^{-1}\mathbf{w} - \frac{1}{2\sigma^{2}}||\mathbf{\Psi}\mathbf{w} - \mathbf{t}||^{2} + \text{const}$$

$$= -\frac{1}{2}\mathbf{w}^{\mathsf{T}}\mathbf{S}^{-1}\mathbf{w} - \frac{1}{2\sigma^{2}}\left(\mathbf{w}^{\mathsf{T}}\mathbf{\Psi}^{\mathsf{T}}\mathbf{\Psi}\mathbf{w} - 2\mathbf{t}^{\mathsf{T}}\mathbf{\Psi}\mathbf{w} + \mathbf{t}^{\mathsf{T}}\mathbf{t}\right) + \text{const}$$

$$= -\frac{1}{2}\mathbf{w}^{\mathsf{T}}\left(\sigma^{-2}\mathbf{\Psi}^{\mathsf{T}}\mathbf{\Psi} + \mathbf{S}^{-1}\right)\mathbf{w} - \frac{1}{\sigma^{2}}\mathbf{t}^{\mathsf{T}}\mathbf{\Psi}\mathbf{w} + \text{const} \text{ (complete the square!)}$$

Thus $\mathbf{w} \mid \mathcal{D} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where

$$\boldsymbol{\mu} = \sigma^{-2} \boldsymbol{\Sigma} \boldsymbol{\Psi}^{\top} \mathbf{t}$$
$$\boldsymbol{\Sigma} = \left(\sigma^{-2} \boldsymbol{\Psi}^{\top} \boldsymbol{\Psi} + \mathbf{S}^{-1} \right)^{-1}$$

- Since a Gaussian prior leads to a Gaussian posterior, this means the Gaussian distribution is the conjugate prior for linear regression!
- Compare μ the closed-form solution for linear regression:

$$\mathbf{w} = \left(\mathbf{\Psi}^{\top}\mathbf{\Psi} + \lambda\mathbf{I}\right)^{-1}\mathbf{\Psi}^{\top}\mathbf{t}$$

This is the mean of the posterior, assuming that $\mathbf{S} = \lambda^{-1} \mathbf{I}$ and $\sigma = 1$. • λ^{-1} is the standard deviation of the prior. As this becomes infinite, the mean of the posterior converges to the maximum likelihood solution.



- Bishop, Pattern Recognition and Machine Learning

• Example with radial basis function (RBF) features

$$\psi_j(x) = \exp\left(-\frac{(x-\mu_j)^2}{2s^2}\right)$$



- Bishop, Pattern Recognition and Machine Learning

Functions sampled from the posterior:



STA314-Lec11

- The posterior just gives us distribution over the parameter space, but if we want to make predictions, the natural choice is to use the posterior predictive distribution.
- Posterior predictive distribution:

$$p(t \mid \mathbf{x}, \mathcal{D}) = \int \underbrace{p(t \mid \mathbf{x}, \mathbf{w})}_{\mathcal{N}(t; \mathbf{w}^{\top} \psi(\mathbf{x}), \sigma)} \underbrace{p(\mathbf{w} \mid \mathcal{D})}_{\mathcal{N}(\mathbf{w}; \boldsymbol{\mu}, \boldsymbol{\Sigma})} d\mathbf{w}$$

Another interpretation: t = w^Tψ(x) + ε, where ε ~ N(0, σ) is independent of w | D ~ N(μ, Σ).

- Another interpretation: t = w^Tψ(x) + ε, where ε ~ N(0, σ) is independent of w | D ~ N(μ, Σ).
- By the linear combination rules for Gaussian random variables, *t* is a Gaussian distribution with parameters

$$\mu_{\text{pred}} = \boldsymbol{\mu}^{\top} \boldsymbol{\psi}(\mathbf{x})$$

$$\sigma_{\text{pred}}^{2} = \boldsymbol{\psi}(\mathbf{x})^{\top} \boldsymbol{\Sigma} \boldsymbol{\psi}(\mathbf{x}) + \sigma^{2}$$

• Hence, the posterior predictive distribution is $\mathcal{N}(t; \mu_{\text{pred}}, \sigma_{\text{pred}}^2)$.

Here we visualize confidence intervals based on the posterior predictive mean and variance at each point:



- The formulation of PCA that we saw earlier in the course was motivated heuristically.
- We will show that it can be expressed as the maximum likelihood estimate of a certain probabilistic model.

Recall: PCA

- Data set $\{\mathbf{x}^{(i)}\}_{i=1}^{N}$
- Each input vector $\mathbf{x}^{(i)} \in \mathbb{R}^{D}$ is approximated as $\hat{\boldsymbol{\mu}} + \mathbf{U}\mathbf{z}^{(i)}$,

$$\mathbf{x}^{(i)} \approx \tilde{\mathbf{x}}^{(i)} = \hat{\boldsymbol{\mu}} + \mathbf{U}\mathbf{z}^{(i)}$$

where $\hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i} \mathbf{x}^{(i)}$ is the data mean, $\mathbf{U} \in \mathbb{R}^{D \times K}$ is the orthogonal basis for the principal subspace, and $\mathbf{z}^{(i)} \in \mathbb{R}^{K}$ is the code vector

$$\mathbf{z}^{(i)} = \mathbf{U}^{\top} (\mathbf{x}^{(i)} - \hat{\mu})$$

• U is chosen to minimize the reconstruction error

$$\mathbf{U}^* = \arg\min_{\mathbf{U}} \sum_{i} ||\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}} + \mathbf{U}\mathbf{U}^{\top}(\mathbf{x}^{(i)} - \hat{\boldsymbol{\mu}})||^2$$

- To formulate probabilistic PCA, let's start with a latent variable model.
- Similar to the Gaussian mixture model, but we will assume continuous, Gaussian latents:

$$\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

 $\mathbf{x} \mid \mathbf{z} \sim \mathcal{N}(\mathbf{W}\mathbf{z} + \boldsymbol{\mu}, \sigma^2 \mathbf{I})$

- Note: this is a naive Bayes model, because $p(\mathbf{x} \mid \mathbf{z})$ factorizes with respect to the dimensions of \mathbf{x} .
- What sort of data does this model produce?

- z is a random coordinate within the affine space centered at μ and spanned by the columbs of **W**.
- To get the random variable x, we samples a standard Normal z and then add a small amount of isotropic noise to Wz + μ.

Probabilistic PCA



- Bishop, Pattern Recognition and Machine Learning

• To perform maximum likelihood in this model, we need to maximize the following:

$$\max_{\mathbf{W},\boldsymbol{\mu},\sigma^{2}} \log p(\mathbf{x} \mid \mathbf{W},\boldsymbol{\mu},\sigma^{2}) = \max_{\mathbf{W},\boldsymbol{\mu},\sigma^{2}} \log \int p(\mathbf{x} \mid \mathbf{z},\mathbf{W},\boldsymbol{\mu},\sigma^{2}) p(\mathbf{z}) d\mathbf{z}$$

- This was hard for the Gaussian mixture model, but in this case it's easy.
- p(x | W, μ, σ²) will be Gaussian (confirm this) and so we just need to compute and Cov[x] and E[x].

$$\mathbb{E}[\mathbf{x}] = \mathbb{E}[\mathbf{W}\mathbf{z} + \boldsymbol{\mu} + \boldsymbol{\epsilon}] = \boldsymbol{\mu}$$

$$Cov[\mathbf{x}] = \mathbb{E}[(\mathbf{W}\mathbf{z} + \epsilon)(\mathbf{W}\mathbf{z} + \epsilon)^{\top}]$$
$$= \mathbb{E}[(\mathbf{W}\mathbf{z}\mathbf{z}^{\top}\mathbf{W}^{\top}] + Cov[\epsilon\epsilon^{\top}]$$
$$= \mathbf{W}\mathbf{W}^{\top} + \sigma^{2}\mathbf{I}$$

• Thus, the likelihood of the data under this model is given by

$$-\frac{ND}{2}\log(2\pi) - \frac{N}{2}\log|\mathbf{C}| - \frac{1}{2}\sum_{i=1}^{N}(\mathbf{x}^{(i)} - \mu)^{\top}\mathbf{C}^{-1}(\mathbf{x}^{(i)} - \mu)$$

where $\mathbf{C} = \mathbf{W}\mathbf{W}^{\mathsf{T}} + \sigma^2 \mathbf{I}$.

• It's a bit involved to derive the maximum likelihood solution, so we will skip it, but Tipping and Bishop (Probabilistic PCA, 1999) show that this is maximized at the following stationary points.

Probabilistic PCA : Maximum Likelihood

$$\hat{\boldsymbol{\mu}}_{\mathsf{MLE}} = \frac{1}{N} \sum_{i=1}^{N} \mathbf{x}^{(i)}$$

$$\hat{\mathbf{W}}_{\mathsf{MLE}} = \hat{\mathbf{U}}_{\mathsf{MLE}} (\hat{\mathbf{L}}_{\mathsf{MLE}} - \hat{\sigma}_{\mathsf{MLE}}^2 \mathbf{I})^{\frac{1}{2}} \mathbf{R}$$

where $\hat{\mathbf{U}}_{MLE}$ is the matrix whose columns are the K unit eigenvectors of the empirical covariance matrix $\hat{\mathbf{\Sigma}}$ that have the largest eigenvalues, $\hat{\mathbf{L}}_{MLE} \in \mathbb{R}^{K \times K}$ is the diagonal matrix whose elements are the corresponding eigenvalues, and \mathbf{R} is any orthogonal matrix.

$$\hat{\sigma}_{\mathsf{MLE}}^2 = \frac{1}{D-K} \sum_{i=K+1}^{D} \lambda_i$$

where λ_i is the *i*th largest eigenvalue of the empirical covariance matrix $\hat{\Sigma}$ of the data. In otherwords, the average variance of the discarded subspace.

- That seems complex, to get an intuition about how this model behaves when it is fit to data, lets consider the MLE density.
- Recall that the marginal distribution on **x** in our fitted model is a Gaussian with mean

$\hat{\mu}_{\mathsf{MLE}}$

and covariance

$$\hat{\mathbf{W}}_{\mathsf{MLE}}\hat{\mathbf{W}}_{\mathsf{MLE}}^{\top} + \hat{\sigma}_{\mathsf{MLE}}^{2}\mathbf{I} = \hat{\mathbf{U}}_{\mathsf{MLE}}(\hat{\mathbf{L}}_{\mathsf{MLE}} - \hat{\sigma}_{\mathsf{MLE}}^{2}\mathbf{I})\hat{\mathbf{U}}_{\mathsf{MLE}}^{\top} + \hat{\sigma}_{\mathsf{MLE}}^{2}\mathbf{I}$$

• The covariance gives us a nice intuition about the type of model this forms.

Probabilistic PCA : Maximum Likelihood

 Consider centering the data and checking the variance along one of the unit eigenvectors u_i, which are the eigenvectors forming the columns of Û_{MLE}:

$$\begin{aligned} \mathsf{Var}(\mathbf{u}_{i}^{\top}(\mathbf{x} - \hat{\boldsymbol{\mu}}_{\mathsf{MLE}})) &= \mathbf{u}_{i}^{\top} \mathsf{Cov}[\mathbf{x}] \mathbf{u}_{i} \\ &= \mathbf{u}_{i}^{\top} \hat{\mathbf{U}}_{\mathsf{MLE}} (\hat{\mathbf{L}}_{\mathsf{MLE}} - \hat{\sigma}_{\mathsf{MLE}}^{2} \mathbf{I}) \hat{\mathbf{U}}_{\mathsf{MLE}}^{\top} \mathbf{u}_{i} + \hat{\sigma}_{\mathsf{MLE}}^{2} \\ &= \lambda_{i} - \hat{\sigma}_{\mathsf{MLE}}^{2} + \hat{\sigma}_{\mathsf{MLE}}^{2} = \lambda_{i} \end{aligned}$$

• Now, consider centering the data and checking the variance along any unit vector orthogonal to the subspace spanned by \hat{U}_{MLE} :

$$Var(\mathbf{u}_{i}^{\top}(\mathbf{x} - \hat{\boldsymbol{\mu}}_{\mathsf{MLE}})) = \mathbf{u}_{i}^{\top}\hat{\mathbf{U}}_{\mathsf{MLE}}(\hat{\mathbf{L}}_{\mathsf{MLE}} - \hat{\sigma}_{\mathsf{MLE}}^{2}\mathbf{I})\hat{\mathbf{U}}_{\mathsf{MLE}}^{\top}\mathbf{u}_{i} + \hat{\sigma}_{\mathsf{MLE}}^{2}$$
$$= \hat{\sigma}_{\mathsf{MLE}}^{2}$$

• In other words, the model captures the variance along the principle axes and approximates the variance in all remaining directions with a single variance.

Intro ML (UofT)

Probably easier to visualize after implementing it.

• The posterior mean is given by

$$\mathbb{E}[\mathbf{z} \,|\, \mathbf{x}] = \left(\hat{\mathbf{W}}_{\mathsf{MLE}}^{\mathsf{T}} \hat{\mathbf{W}}_{\mathsf{MLE}} + \hat{\sigma}_{\mathsf{MLE}}^{2} \mathbf{I} \right)^{-1} \hat{\mathbf{W}}_{\mathsf{MLE}}^{\mathsf{T}} (\mathbf{x} - \hat{\mu}_{\mathsf{MLE}})$$

 $\bullet\,$ So, if we don't fit σ^2 and instead take it to 0 we get

$$\mathbb{E}[\mathbf{z} \mid \mathbf{x}] \stackrel{\sigma^2 \to 0}{\to} \left(\hat{\mathbf{W}}_{\mathsf{MLE}}^{\mathsf{T}} \hat{\mathbf{W}}_{\mathsf{MLE}} \right)^{-1} \hat{\mathbf{W}}_{\mathsf{MLE}}^{\mathsf{T}} (\mathbf{x} - \hat{\mu}_{\mathsf{MLE}})$$

 $\bullet\,$ Can show that this is a projection onto an affine space spanned by the columns of $\hat{U}_{MLE}.$

- Fitting a full-covariance Gaussian model of data requires D(D + 1)/2 + D parameters. With PPCA we model only the K most significant correlations and this only requires O(D) parameters as long as K is small.
- Basis of Bayesian treatement of PCA, which gives us a Bayesian method for determining the dimensionality of the principal subspace (i.e. K).
- Existence of likelihood functions allows direct comparison with other probabilistic models.
- Can use PPCA as a class-conditional density (as in GDA) to reduce the requirement to fit and store $\mathcal{O}(D^2)$ parameters.

- Gaussian discriminant analysis.
 - ► Gaussian class-conditional generative model p(x | t) used for classification.
- Gaussian mixture model.
 - Gaussian latent variable model $p(\mathbf{x}) = \sum_{z} p(\mathbf{x}, z)$ used for clustering.
- Bayesian linear regression.
 - Gaussian discriminative model p(t | x) used for regression with a Bayesian analysis for the weights.
- Probabilistic PCA.
 - Gaussian latent variable model p(x) = ∫_z p(x, z) used for dimensionality reduction.

Optional material: Bayesian model selection

- Consider selecting models from a Bayesian perspective.
- Related to Occam's Razor: "Entities should not be multiplied beyond necessity."
 - Named after the 14th century British theologian William of Occam
- Huge number of attempts to formalize mathematically
 - See Domingos, 1999, "The role of Occam's Razor in knowledge discovery" for a skeptical overview.

https://homes.cs.washington.edu/~pedrod/papers/dmkd99.pdf

- Common misinterpretation: your prior should favor simple explanations
- Better interpretation: by averaging over many hypothesis, Bayesian model selection naturally prefers simpler models.

- Suppose you have a finite set of models, or hypotheses {\$\mathcal{H}_i\$}_{i=1}^M\$ (e.g. polynomials of different degrees)
- Posterior inference over models (Bayes' Rule):

$$p(\mathcal{H}_i \mid \mathcal{D}) \propto \underbrace{p(\mathcal{H}_i)}_{\text{prior}} \underbrace{p(\mathcal{D} \mid \mathcal{H}_i)}_{\text{evidence}}$$

• The evidence is also called marginal likelihood since it requires marginalizing out the parameters:

$$p(\mathcal{D} \mid \mathcal{H}_i) = \int p(\mathbf{w} \mid \mathcal{H}_i) p(\mathcal{D} \mid \mathbf{w}, \mathcal{H}_i) \, \mathrm{d}\mathbf{w}$$

- p(H_i) is typically uniform, so we can compare them based on marginal likelihood.
- Bayesian model selection:

$$\mathcal{H}^* = \arg\max_i p(\mathcal{D} \mid \mathcal{H}_i)$$

• What types of models does this procedure prefer?

Occam's Razor (optional)

- Suppose M_1 , M_2 , and M_3 denote a linear, quadratic, and cubic model.
- M_3 is capable of explaning more datasets than M_1 .
- But its distribution over \mathcal{D} must integrate to 1, so it must assign lower probability to ones it can explain.



⁻ Bishop, Pattern Recognition and Machine Learning

Occam's Razor (optional)

• How does the evidence penalize complex models?



• Approximating the integral for $\mathbf{w} \in \mathbb{R}$:

$$p(\mathcal{D} | \mathcal{H}_i) = \int p(\mathcal{D} | \mathbf{w}, \mathcal{H}_i) p(\mathbf{w} | \mathcal{H}_i)$$

$$\approx \underbrace{p(\mathcal{D} | \mathbf{w}_{MAP}, \mathcal{H}_i)}_{\text{best-fit likelihood}} \underbrace{\frac{\Delta \mathbf{w}_{\text{posterior}}}{\Delta \mathbf{w}_{\text{prior}}}}_{\text{Occam factor}}$$

Let's investigate

$$\log p(\mathcal{D} \mid \mathcal{H}_i) = \log p(\mathcal{D} \mid \mathbf{w}_{MAP}, \mathcal{H}_i) + \log \frac{\Delta \mathbf{w}_{posterior}}{\Delta \mathbf{w}_{prior}}$$

- First term represents fit to the data given the most probable parameter values.
- Second term is a penalty, because it is negative (Δw_{posterior} < Δw_{prior}).
- Thus if the posterior is very peaked and confident about the data, this penalty term will be very negative.

• $\mathbf{w} \in \mathbb{R}^{M}$ we have

$$\log p(\mathcal{D} \mid \mathcal{H}_i) = \log p(\mathcal{D} \mid \mathbf{w}_{MAP}, \mathcal{H}_i) + M \log \frac{\Delta \mathbf{w}_{\text{posterior}}}{\Delta \mathbf{w}_{\text{prior}}}$$

- So the more parameters we have, the higher the penalty.
- Optimal model complexity is determined by a tradeoff.
- In Bayesian model selection, we naturally prefer simpler models that model the data well.