

MLPR Week 6a

Bayesian linear regression **prior**:

draw weight, bias from gaussians

= **gaussian process** (defining function value) with specific kernel specified by parameters σ_w^2 and σ_b^2 (hyperparameters)

weight view ($\mathbf{W}^T \mathbf{X}$) & function value view (f_i)

$$\tilde{f}_i = f(\mathbf{x}^{(i)}) = \mathbf{w}^T \mathbf{x}^{(i)} + b, \quad \mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma_w^2 \mathbf{I}), \quad b \sim \mathcal{N}(0, \sigma_b^2)$$

$$\begin{aligned} \text{cov}(\tilde{f}_i, \tilde{f}_j) &= \mathbb{E}[\tilde{f}_i \tilde{f}_j] - \overbrace{\mathbb{E}[\tilde{f}_i] \mathbb{E}[\tilde{f}_j]}^{0 \dots} \\ &= \mathbb{E}[(\mathbf{w}^T \mathbf{x}^{(i)} + b)^T (\mathbf{w}^T \mathbf{x}^{(j)} + b)] \\ &= \sigma_w^2 \mathbf{x}^{(i)T} \mathbf{x}^{(j)} + \sigma_b^2 = k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}). \end{aligned}$$

Since

Line 3: Since expectation over w , could take out x
also line 3: terms in " \dots " are combination of w and b ,
since w and b are independent, " \dots " = 0

$$\begin{aligned} &= \mathbb{E}[(\mathbf{w}^T \mathbf{x}^{(i)} + b)^T (\mathbf{w}^T \mathbf{x}^{(j)} + b)] \\ &= \mathbb{E}[\mathbf{x}^{(i)T} \mathbf{w} \mathbf{w}^T \mathbf{x}^{(j)} + b^2 + \dots] \\ &= \mathbf{x}^{(i)T} \underbrace{\mathbb{E}[\mathbf{w} \mathbf{w}^T]}_{\sigma_w^2 \mathbf{I}} \mathbf{x}^{(j)} + \underbrace{\mathbb{E}[b^2]}_{\sigma_b^2} + \underbrace{\dots}_0 \\ \Rightarrow k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) &= \sigma_w^2 \underbrace{\mathbf{x}^{(i)T} \mathbf{x}^{(j)}}_{\phi(\mathbf{x}^{(i)})^T \phi(\mathbf{x}^{(j)})} + \sigma_b^2 \end{aligned}$$

kernel trick

Drawback of finite RBFs in Bay linear regression:

1. selected manually, not good for exploration
2. If x move away from origin, always zero ("under fitting")

put RBFs everywhere

= **infinite RBFs**, = **mercer kernel**

infinite dimension feature vector, but **take inner product**, only get scalar value

=> kernel trick / kernelized: Replacing an inner product with a kernel function

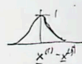
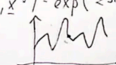
IF we put RBFs everywhere
 ~> Analytically derive $\underline{\phi}^T \underline{\phi}$
 $k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \sigma^2 \underline{\phi}(\underline{x}^{(i)})^T \underline{\phi}(\underline{x}^{(j)})$
 $\downarrow \# \text{ RBFs} \rightarrow \infty$
 $k(\underline{x}^{(i)}, \underline{x}^{(j)}) \propto \exp(-\|\underline{x}^{(i)} - \underline{x}^{(j)}\|^2)$
 "kernel trick"
 - Rewrite your algo. so that it only depends on inner products
 near Set $\underline{\phi}(\underline{x}^{(i)})^T \underline{\phi}(\underline{x}^{(j)}) = \overbrace{k(\underline{x}^{(i)}, \underline{x}^{(j)})}^{\text{Mercer kernel}}$

GP Prediction: in week 5 bottom

Can we be more uncertain at prediction in response to a surprising training label?
 => change GP kernel function, learn the parameters in kernel

Choosing a family of kernel functions

Kernel function examples

- Squared exponential
 $k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \exp(-\|\underline{x}^{(i)} - \underline{x}^{(j)}\|^2)$

- Periodic kernel: *don't minimize*
 $k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \exp(-2 \sin(\pi \|\underline{x}^{(i)} - \underline{x}^{(j)}\| / \tau))$

- Kernels can be combined:
 $k(\underline{x}^{(i)}, \underline{x}^{(j)}) = \alpha k_1(\underline{x}^{(i)}, \underline{x}^{(j)}) + \beta k_2(\underline{x}^{(i)}, \underline{x}^{(j)})$
 $\alpha, \beta > 0$
 is kernel function if k_1, k_2 kernel function

Another kernel function with learnable parameters: RBF kernel
hyperparameters σ_f (marginal variance), ℓ_d (curve width)

$$k(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = \sigma_f^2 \exp\left(-\frac{1}{2} \sum_{d=1}^D (x_d^{(i)} - x_d^{(j)})^2 / \ell_d^2\right)$$

The **marginal variance** of **one function value** \tilde{f}_i is:

$$\text{var}[\tilde{f}_i] = k(\mathbf{x}^{(i)}, \mathbf{x}^{(i)}) = \sigma_f^2.$$

larger marginal variance = larger width in y dimension
smaller l_d , get more wiggly curve

Learning the hyperparameters: choose the parameters (σ_f and l_d) with largest marginal likelihood (calculate marginal likelihood = just copy values)

Pick parameters by (marginal) likelihood:

$$p(y|X, \theta = \{\sigma_y^2, \sigma_f^2, \{l_d\}\})$$
$$= \mathcal{N}(y; \underline{0}, K(X, X) + \sigma_y^2 \mathbf{I})$$

so $\log P(y | x, M) = \log$ of the standard multivariate Gaussian pdf

covariance = kernel matrix evaluated at the training inputs + observation noise

fully Bayesian approach integrating over all possible hyperparameters, but can't be computed exactly

Week6b

gradient descent not constrained,

If step size not small enough, could take weight into illegal value / not converge

=> if got constrained parameters (e.g. positive), take log (exponential)
obtain that derivative with the chain rule

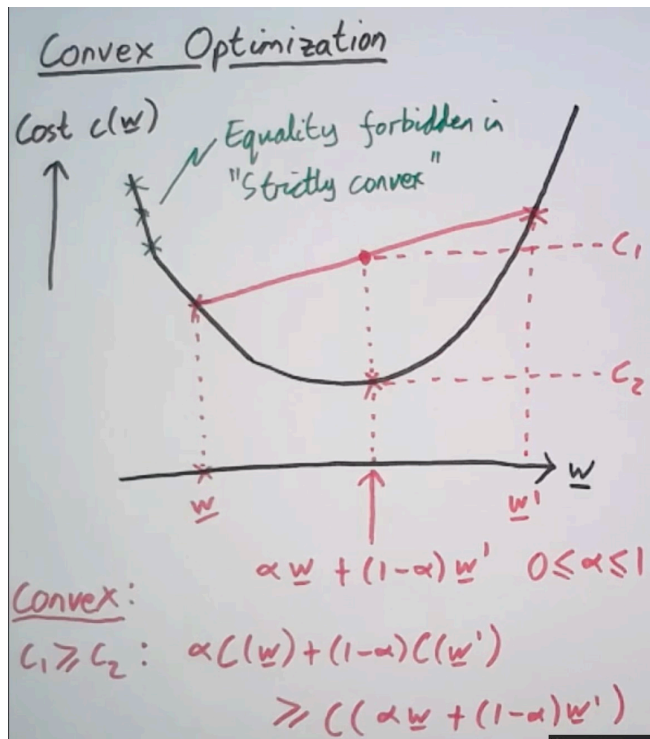
Convex:

if convex, guarantee minimum

convex (could have straight line) VS strictly convex

norms are convex

sum of any convex functions is convex



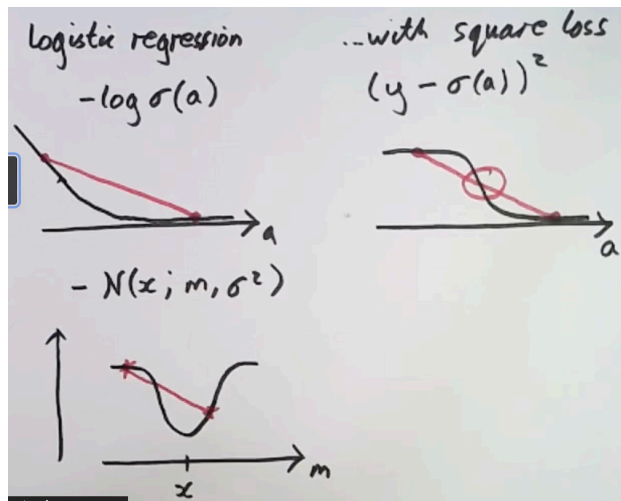
Not convex:

Square loss, not convex optimization problem, harder to optimize

why take log instead of just optimize probability:

gaussian is not convex (one reason)

++ (sum over examples, avoid numerical underflow)



if not convex, could still try to reduce this loss function using gradient methods

Week6c How to develop Softmax loss function (regression):

Define likelihood $P(Y | X, W)$, want positive score, normalized to sum to 1

log-probability of a single observation:

C: y of the current sample

K: the current class when doing sum

Kronecker delta δ_{kc} : if k and c the same, get 1

finally take x out

$$\begin{aligned}\log P(y_c=1 | \mathbf{x}, W) &= \log f_c = (\mathbf{w}^{(c)})^\top \mathbf{x} - \log \sum_{k'} e^{(\mathbf{w}^{(k')})^\top \mathbf{x}} \\ \nabla_{\mathbf{w}^{(k)}} \log f_c &= \delta_{kc} \mathbf{x} - \frac{1}{\sum_{k'} e^{(\mathbf{w}^{(k')})^\top \mathbf{x}}} \mathbf{x} e^{(\mathbf{w}^{(k)})^\top \mathbf{x}} \\ &= (y_k - f_k) \mathbf{x}.\end{aligned}$$

Derivative: if label & function output different, move in the direction of the input
decrease the scorer of all other classes (in denominator of softmax function)

++ dimension of w: each row is $\mathbf{w}^{(k)}$, each col is feature, total $K \times D$

Redundant parameters: (when divided by numerator)

when only have 2 classes, = logistic regression

but only depends on difference between $\mathbf{w}^{(1)}$ and $\mathbf{w}^{(2)}$, could set $\mathbf{w}^{(2)} = 0$

in multi class case, weight matrix is $K \times D$, could set one row = 0

in neural network most parameters are redundant (over-parameterized)

when have over-parameterized model, don't try to interpret parameters

$$\begin{aligned}P(y=1 | \mathbf{x}, W) &= \frac{e^{(\mathbf{w}^{(1)})^\top \mathbf{x}}}{e^{(\mathbf{w}^{(1)})^\top \mathbf{x}} + e^{(\mathbf{w}^{(2)})^\top \mathbf{x}}} \\ &= \frac{1}{1 + e^{(\mathbf{w}^{(2)} - \mathbf{w}^{(1)})^\top \mathbf{x}}} = \sigma((\mathbf{w}^{(1)} - \mathbf{w}^{(2)})^\top \mathbf{x}).\end{aligned}$$

Week6d Robust logistic regression

have corrupted label, want to describe this mathematically, derive a different loss

binary choice m ($m=1$, correct / $m=0$, corrupt label) for each observation
 since true label hidden

epsilon = very small probability

Assume Epsilon is entirely random, make the choice independent of the input position and the weights (or...)

$$P(m | \epsilon) = \text{Bernoulli}(m; 1-\epsilon) = \begin{cases} 1-\epsilon & m = 1 \\ \epsilon & m = 0. \end{cases}$$

If correct label ($m=1$), do nothing, if corrupt label, output at random (Or...)

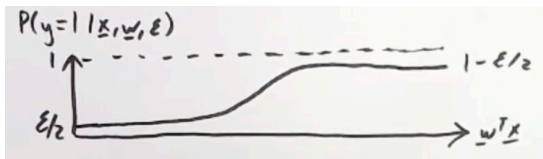
$$P(y=1 | \mathbf{x}, \mathbf{w}, m) = \begin{cases} \sigma(\mathbf{w}^\top \mathbf{x}) & m = 1 \\ \frac{1}{2} & m = 0. \end{cases}$$

Introduce m by sum & product rule

After product rule, cancel out \mathbf{x} & \mathbf{w} , since m does not depends on those

$$\begin{aligned} P(y=1 | \mathbf{x}, \mathbf{w}, \epsilon) &= \sum_{m \in \{0,1\}} P(y=1, m | \mathbf{x}, \mathbf{w}, \epsilon) \\ &= \sum_{m \in \{0,1\}} P(y=1 | \mathbf{x}, \mathbf{w}, m) P(m | \epsilon) \\ &= \underbrace{(1-\epsilon)}_{m=1} \underbrace{\sigma(\mathbf{w}^\top \mathbf{x})}_{m=0} + \epsilon \frac{1}{2}. \end{aligned}$$

Result: model could not be 100% confident in prediction



Get the Gradient, = correction term + normal logistic gradient

Tweak epsilon

get epsilon

