

Week 3

Bayes & Naive Bayes

Week 3b Regression and Gradients

Residual & minimize sum of squared residuals

$$\begin{aligned}\mathbf{r}^\top \mathbf{r} &= (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) \\ &= \mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X}\mathbf{w}.\end{aligned}$$

Find the gradient: (partial derivative, wT gone)

$$\nabla_{\mathbf{w}} [\mathbf{r}^\top \mathbf{r}] = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\mathbf{w}.$$

(why delta W [WT h] = h:)

$$\begin{aligned}\nabla_{\mathbf{w}} \left[\frac{\partial \mathbf{r}^\top \mathbf{r}}{\partial \mathbf{w}} \right] &= \begin{bmatrix} \frac{\partial \mathbf{r}^\top \mathbf{r}}{\partial w_1} \\ \frac{\partial \mathbf{r}^\top \mathbf{r}}{\partial w_2} \\ \vdots \\ \frac{\partial \mathbf{r}^\top \mathbf{r}}{\partial w_n} \end{bmatrix} = \underline{\mathbf{h}} \\ \frac{\partial \mathbf{r}^\top \mathbf{r}}{\partial w_i} &= \frac{\partial}{\partial w_i} \sum_j w_j h_j = \sum_j \underline{w_i} h_j \\ &= \frac{\partial}{\partial w_i} (w_1 h_1 + w_2 h_2 + \dots + w_i h_i + \dots + w_n h_n) \\ &= h_i\end{aligned}$$

set all gradient = 0 (move in any direction can't improve the cost), get closed form solution

Assume inverse $(\mathbf{X}^\top \mathbf{X})^{-1}$ exist (able to compute it)

Assume only one weight could solve this equation

need $N \geq D$, and datapoints not linear dependent

since if $N < D$ (or D linear dependent), then matrix $(\mathbf{X}^\top \mathbf{X})^{-1}$ is low rank (rank N), not invertible

solution: L2 regularization (add datapoints)

$$\begin{aligned}-2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X}\mathbf{w} &\stackrel{\text{...}}{=} \mathbf{0} \\ \Rightarrow \quad \mathbf{X}^\top \mathbf{X}\mathbf{w} &= \mathbf{X}^\top \mathbf{y} \\ \Rightarrow \quad \mathbf{w} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}.\end{aligned}$$

why can't expand the inverse and cancel out terms, like this:

$$\begin{aligned}&= \cancel{\mathbf{X}^{-1} \mathbf{X}^\top} \cancel{\mathbf{X}^\top} \mathbf{y} \\ &= \cancel{\mathbf{X}^{-1} \mathbf{y}} ? \\ &\quad \cancel{\mathbf{X}} \quad \text{NxD}\end{aligned}$$

since X is N by D matrix, usually we have many more datapoints than dimension, can't invert (inverse only exist for square matrix)

$(X^T X)^{-1} X^T$ is pseudo inverse

steepest-descent method:

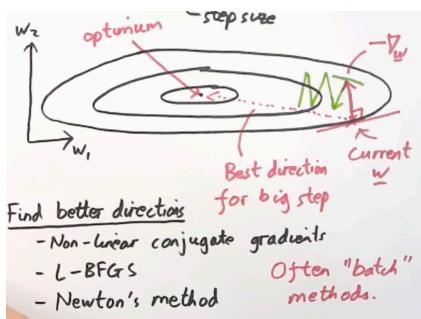
$$\mathbf{w} \leftarrow \mathbf{w} - \eta \nabla_{\mathbf{w}} [\mathbf{r}^T \mathbf{r}],$$

Drawback: Gradient descent not in sensible direction (zig zag)

Why gradient descent instead of closed form solution:

Numerical algorithm takes long time, have accuracy guarantees

sometimes, iterative algorithms can have good tradeoffs in accuracy & time (if want guarantee "accurately solving most difficult numerical problem")



rewrite the gradient: (batch gradient descent)

$$\begin{aligned}\nabla_{\mathbf{w}} [\mathbf{r}^T \mathbf{r}] &= 2X^T(X\mathbf{w}) - 2X^T\mathbf{y} \\ &= 2X^T(\mathbf{f} - \mathbf{y}) \\ &= 2 \sum_n \mathbf{x}^{(n)} (f^{(n)} - y^{(n)}).\end{aligned}$$

Mini batch gradient descent: mini batch size B (SGD: $B = 1$)

since batch GD expensive: one step need to go through all points

$$\begin{aligned}&= -2N \frac{1}{N} \sum_n \mathbf{x}^{(n)} (y^{(n)} - f^{(n)}) \\ &\approx -2N \frac{1}{B} \sum_{b=1}^B \mathbf{x}^{(b)} (y^{(b)} - f^{(b)})\end{aligned}$$

weights are pulled most in the direction of inputs that had large prediction errors ($f - y$)

why: look at matrix dimensions & interpretation for mat mul

SGD, each sample gives a crude one-sample Monte Carlo approximation
could drop N, since N could combine into step size alpha

$$\frac{1}{N} \nabla_{\mathbf{w}} [\mathbf{r}^\top \mathbf{r}] \approx 2\mathbf{x}^{(n)}(f^{(n)} - y^{(n)}).$$

Week 3c

Logistic regression, fit binary labels with linear regression
transform output into [0,1]
(could replace x with basis function)

$$f(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w}^\top \mathbf{x}}}.$$

Loss: could use **least square**

with logistic sigmoid, y bounded to [0,1], function values = probability (if got enough data & flexible enough function)

Before, this is not very good interpretation, since if fit RBFs, would get probability <0, or function value >1, can't be probability

Could use square loss again:
 $\sum_{n=1}^N (y^{(n)} - f(\mathbf{x}^{(n)}; \mathbf{w}))^2$

Interpretation:
 $p(y=1|\mathbf{x}) \approx f(\mathbf{x}; \mathbf{w})$

Given this interpretation (function output = probability), more popular to **maximize likelihood** (of the parameters $P(D | w)$), maximize prob of data given w:

$$L(\mathbf{w}) = \prod_{n=1}^N P(y^{(n)} | \mathbf{x}^{(n)}, \mathbf{w}),$$

Least square & maximum likelihood:

different ways to fit weights, got different parameters

interpret one outlier differently: square error doesn't change much with summing an outlier,

but in MLE, multiply big product by number close to 0 changes the whole loss

what if use sum of log likelihood?

Advantage of MLE: if lots of data, converge faster; for small dataset, might prefer square loss

Minimize negative log likelihood:

why sum of log likelihood: could apply Monte Carlo trick to do SGD

Normally write in 2 terms ($y=1$ and $y=0$):

$$\text{NLL} = -\log L(\mathbf{w}) = - \sum_{n=1}^N \log \left[\sigma(\mathbf{w}^\top \mathbf{x}^{(n)})^{y^{(n)}} (1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(n)}))^{1-y^{(n)}} \right],$$

or

$$\text{NLL} = - \sum_{n:y^{(n)}=1} \log \sigma(\mathbf{w}^\top \mathbf{x}^{(n)}) - \sum_{n:y^{(n)}=0} \log(1 - \sigma(\mathbf{w}^\top \mathbf{x}^{(n)})).$$

But a **trick**: here use label $z \{-1, 1\}$, use $z = 2y - 1$ to transform from $y \{0, 1\}$

when label = 1, just sigmoid($\mathbf{W}^\top \mathbf{X}$),

when label = 0, want 1-sigmoid, = sigmoid($-\mathbf{W}^\top \mathbf{X}$)

so could simply put 1/-1 label Z inside

There is a trick to write the cost function more compactly. We transform the **labels to be** $z^{(n)} \in \{-1, +1\}$ where $z^{(n)} = (2y^{(n)} - 1)$, and noticing $\sigma(-a) = 1 - \sigma(a)$, we can write:

$$\text{NLL} = - \sum_{n=1}^N \log \sigma(z^{(n)} \mathbf{w}^\top \mathbf{x}^{(n)}). \quad (6)$$

Gradient: (by chain rule)

$$\begin{aligned} \nabla_{\mathbf{w}} \text{NLL} &= - \sum_{n=1}^N \nabla_{\mathbf{w}} \log \sigma_n = - \sum_{n=1}^N \frac{1}{\sigma_n} \nabla_{\mathbf{w}} \sigma_n = - \sum_{n=1}^N \frac{1}{\sigma_n} \sigma_n (1 - \sigma_n) \nabla_{\mathbf{w}} z^{(n)} \mathbf{w}^\top \mathbf{x}^{(n)}, \\ &= - \sum_{n=1}^N (1 - \sigma_n) z^{(n)} \mathbf{x}^{(n)}. \end{aligned}$$

$$\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a))$$

when classifier is confident and correct on an example, it contributes little to the gradient

SGD improve the examples the classifier gets wrong

Check derivatives https://media.ed.ac.uk/media/1_kyd8kz63

Week 4 Bayesian regression

Bayesian classifier: look for a particular set of parameters, do classification with this set of parameters

Bayesian methods: put beliefs on parameters by probability distribution

Probabilistic model for regression: just a Gaussian around regression line $f(\mathbf{x}, \mathbf{w})$
assume variance sigma y is known and constant

$$p(y | \mathbf{x}, \mathbf{w}) = \mathcal{N}(y; f(\mathbf{x}; \mathbf{w}), \sigma_y^2),$$

MLE, negative log likelihood of normal distribution:

Sum of log likelihood of individual observations: assume observations are independent

assume constant variance, = minimizing square error (linear regression + constant noise)

$$\begin{aligned} -\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) &= -\sum_n \log p(y^{(n)} | \mathbf{x}^{(n)}, \mathbf{w}) \\ &= \frac{1}{2\sigma_y^2} \sum_{n=1}^N [(y^{(n)} - f(\mathbf{x}^{(n)}; \mathbf{w}))^2] + \frac{N}{2} \log(2\pi\sigma_y^2). \end{aligned}$$

Since (log of normal distribution PDF)

$$= \sum_{n=1}^N \underbrace{\frac{1}{2\sigma_y^2} (y^{(n)} - f(\mathbf{x}^{(n)}; \mathbf{w}))^2}_{\text{Term 1}} + \sum_n \underbrace{\frac{1}{2} \log(2\pi\sigma_y^2)}_{\text{Term 2}}$$

if different variance for different datapoint, sigma y is now depends on data point n , can't put this terms outside the sum:

$$(\sigma_y^{(n)})$$

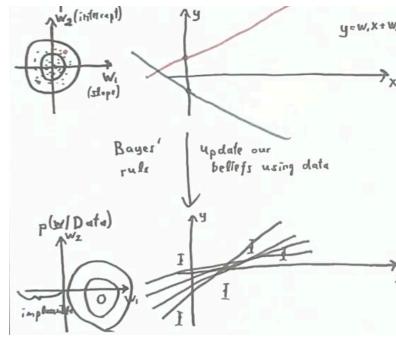
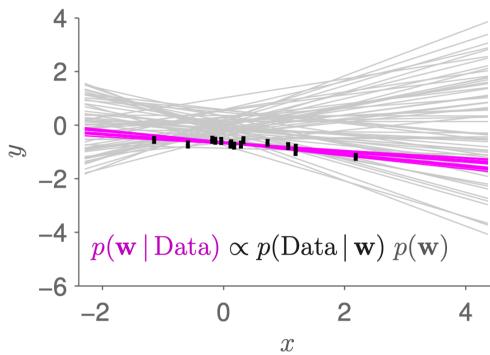
Interpretation: squared difference of each sample weighted by its precision:

if a training point have greater precision, closer to line $f(\mathbf{x}, \mathbf{w})$, constrain the resulting lines a lot more, **more weight it should carry**

Precision:

$$\frac{1}{(\sigma_y^{(n)})^2}$$

systematically express the uncertainty based on limited and noisy data



Bayesian regression:

prior e.g. each point of \mathbf{w} = a function $f(x, \mathbf{w})$

(think of parameter space & data point space)

Where to get prior belief: prior knowledge about the problem / have seen some data beforehand

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, 0.4^2 \mathbb{I}),$$

posterior by bayesian rule, inputs X fixed

likelihood $p(D | \mathbf{w})$ or $p(y | \mathbf{w}, X)$ (= product over each datapoint)

$$p(\mathbf{w} | \mathcal{D}) = p(\mathbf{w} | \mathbf{y}, X) = \frac{p(\mathbf{y} | \mathbf{w}, X) p(\mathbf{w})}{p(\mathbf{y} | X)} \propto p(\mathbf{y} | \mathbf{w}, X) p(\mathbf{w})$$

conjugate prior: product of the prior and likelihood, get same functional form as the prior

get Gaussian posterior:

$$\begin{aligned} p(\mathbf{w} | \mathcal{D}) &\propto p(\mathbf{w}) p(\mathbf{y} | \mathbf{w}, X), \\ &\propto \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2 \mathbb{I}) \prod_n \mathcal{N}(y^{(n)}; \mathbf{w}^\top \boldsymbol{\phi}(x^{(n)}), \sigma_y^2), \\ &\propto \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2 \mathbb{I}) \mathcal{N}(\mathbf{y}; \Phi \mathbf{w}, \sigma_y^2 \mathbb{I}), \quad \text{Vector rep} \end{aligned}$$

$$p(\mathbf{w} | \mathcal{D}) = \mathcal{N}(\mathbf{w}; \mathbf{w}_N, V_N)$$

posterior mean & covariance (gives numerical result)

$$V_N = \sigma_y^2 (\sigma_y^2 V_0^{-1} + \Phi^\top \Phi)^{-1}$$

$$\mathbf{w}_N = V_N V_0^{-1} \mathbf{w}_0 + \frac{1}{\sigma_y^2} V_N \Phi^\top \mathbf{y}.$$

Week 4b, Bayesian inference and prediction how to use the generated distribution

Solve prediction problems: if can't immediately write down expressions for $P(y | \text{data})$,

introduce model parameters and/or latent variables z
by sum rule

$$P(y | \text{data}) = \sum_z P(y, z | \text{data})$$

by product rule (z already known from right probability)

$$P(y | \text{data}) = \sum_z P(y | z, \text{data}) P(z | \text{data})$$

Prediction $P(Y_{\text{test}} | X_{\text{test}}, D)$:

for unknown y in test location x , with training data D , use sum rule & product rule to introduce w (is continuous)

$P(y|x, w)$ instead of $P(y|x, w, D)$; $P(w|D)$ instead of $P(w|x, D)$: doesn't depend on D & test location x

$$p(y | x, D) = \int p(y, w | x, D) dw.$$

$$p(y | x, D) = \int p(y | x, w) p(w | D) dw.$$

here $p(w|D) =$ the Gaussian posterior above

the other term (observation model $p(Y_{\text{test}} | X_{\text{test}}, w)$): gaussian noise around a function $f(x, w)$:

$$p(y | x, w) = \mathcal{N}(y; w^\top x, \sigma_y^2).$$

w could be very high dimension, need to **solve high dimension integral** here

2 ways to solve integral: Mechanistic way & Easy way

Mechanistic way:

inside the integral is the product of 2 gaussian

use tedious linear algebra to find mean & covariance of this **joint distribution**

$$\begin{aligned} p(y | x, D) &= \int p(y, w | x, D) dw \\ &= \dots \text{tedious linear algebra} \\ &= \mathcal{N}\left(\left[\begin{array}{c} w \\ y \end{array}\right]; \left[\begin{array}{cc} m & r \\ r & r^2 \end{array}\right], \left[\begin{array}{cc} \Sigma_{w,w} & \Sigma_{w,y} \\ \Sigma_{y,w} & \Sigma_{y,y} \end{array}\right]\right) dw \\ &= \mathcal{N}(y; m, r^2) \end{aligned}$$

Easy way:

observation model $P(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathbf{w})$ have particular form: regression line $f(\mathbf{x}) +$ constant noise

belief about the function form: $f(\mathbf{x}) = \mathbf{X}^T \mathbf{w}$ (\mathbf{x} is one test example, D dimension)

$$y = f(\mathbf{x}) + \nu = \mathbf{x}^T \mathbf{w} + \nu, \quad \nu \sim \mathcal{N}(0, \sigma_y^2)$$

belief about \mathbf{w} :

$$p(\mathbf{w} | \mathcal{D}) = \mathcal{N}(\mathbf{w}; \mathbf{w}_N, V_N)$$

2 steps to get $P(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D})$:

1. consider belief about function f

f is linear combination of \mathbf{w} , \mathbf{w} is gaussian, so $P(f)$ is gaussian

=> $p(f | \mathbf{X}_{\text{test}}, \mathcal{D})$ is gaussian,

mean is linear combination $\mathbf{x}^T \mathbf{w}_N$

covariance is general equation of how covariance changes given a linear transformation (because...)

Distribution over f at test point \mathbf{x} :

$$p(f | \mathbf{x}, \mathcal{D}) = \mathcal{N}(f; \mathbf{x}^T \mathbf{w}_N, \mathbf{x}^T V_N \mathbf{x})$$

because $\text{Cov}[A\mathbf{w} + b] = A \text{Cov}[\mathbf{w}] A^T$

$$\mu_f = \mathbb{E}[\mathbf{x}^T \mathbf{w}] = \mathbf{x}^T \mathbf{w}_N \quad \text{var}[f] = \mathbf{x}^T V_N \mathbf{x},$$

$$\text{var}[f] = \mathbb{E}[(f - \mu_f)(f - \mu_f)] = \mathbb{E}[\mathbf{x}^T (\mathbf{w} - \mathbf{w}_N)(\mathbf{w} - \mathbf{w}_N)^T \mathbf{x}] = \mathbf{x}^T \text{cov}[\mathbf{w}] \mathbf{x}$$

2. Consider belief about y :

Final belief $P(y_{\text{test}} | \mathbf{x}_{\text{test}}, \mathcal{D}) = p(f | \mathbf{X}_{\text{test}}, \mathcal{D}) + \text{gaussian noise } \nu$

$$p(y | \mathcal{D}, \mathbf{x}) = \mathcal{N}(y; \mathbf{x}^T \mathbf{w}_N, \mathbf{x}^T V_N \mathbf{x} + \sigma_y^2).$$

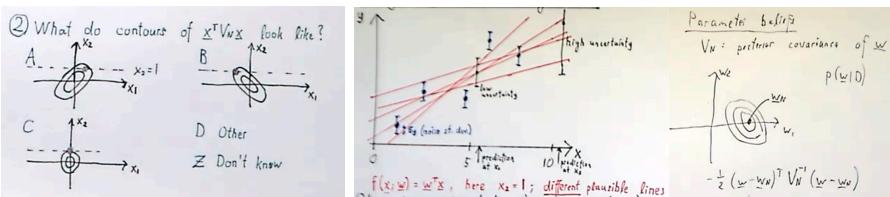
mean, covariance corresponds to m, r^2 in mechanistic way

covariance: allow us to quantify how certain we are about prediction (original goal)

covariance becomes smaller the more data we get, = more certain about the line model; but always have sigma y (when observing infinite data, only sigma y)

larger x, larger covariance, larger uncertainty (most certain at x=0, when bias=0 (according to graph))

find contour of covariance in x space: look at x VS y graph



Make decisions: To have one output y (instead of just Bayesian approach):

Loss function: $L(y, \hat{y})$ (\hat{y} is point estimate)

Cost function: expected loss given training data D

$$c = \mathbb{E}_{p(y|x, \mathcal{D})}[L(y, \hat{y})] = \int L(y, \hat{y}) p(y | x, \mathcal{D}) dy.$$

if squared loss, differentiate the cost w.r.t y

$$L(y, \hat{y}) = (y - \hat{y})^2$$

$$\frac{\partial c}{\partial \hat{y}} = \mathbb{E}_{p(y|x, \mathcal{D})}[-2(y - \hat{y})] = -2(\mathbb{E}_{p(y|x, \mathcal{D})}[y] - \hat{y})$$

set derivative = 0, **point estimate \hat{y} = posterior mean**

$$\hat{y} = \mathbb{E}_{p(y|x, \mathcal{D})}[y]$$

posterior mean corresponds to using an L2 regularized model (but uncertainty of bayesian still useful)

(MSE: symmetric loss, if underestimate the true value, cost would be as big as if you overestimate it, might not appropriate)

(e.g. when the cost of underestimating sales is much higher than the cost of overestimating)