

Linear Models II

CS771: Introduction to Machine Learning

Purushottam Kar

Recap of Last Lecture

2

Looked at SGD as a speedier way to implement GD

Mini-batch SGD as a more stable version of SGD

Dealing with constraints in optimization problems

Projected Gradient Descent: simple but only if projection is easy

Lagrangian Dual Problems: notions of dual variables/Lagrange multipliers, Lagrangian, primal/dual problems

Strong duality, complementary slackness

Looked at the dual problem for the CSVM (with/without bias/slack)



CSVM Dual Problem

3

If we have a bias b , then the dual problem looks like

$$\begin{aligned} & \underset{\boldsymbol{\alpha} \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \boldsymbol{\alpha}_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \boldsymbol{\alpha}_i \boldsymbol{\alpha}_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \right\} \\ & \text{s.t. } \boldsymbol{\alpha}_i \in [0, C], \text{ and } \sum_{i=1}^n \boldsymbol{\alpha}_i y^i = 0 \end{aligned}$$

If we omit bias (or hide it inside the model vector itself) dual is

$$\begin{aligned} & \underset{\boldsymbol{\alpha} \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \boldsymbol{\alpha}_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \boldsymbol{\alpha}_i \boldsymbol{\alpha}_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \right\} \\ & \text{s.t. } \boldsymbol{\alpha}_i \in [0, C] \end{aligned}$$

Let us see a method to solve this problem



Coordinate Descent

4

Similar to GD except only one coordinate is changed in a single step

For example, let us take the problem $\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$ s.t. $\mathbf{x} \in \mathcal{C}$

Let $\nabla_j f(\mathbf{x}) = \frac{\partial f}{\partial \mathbf{x}_j}$ be the j th partial derivative

CCD: choose coordinate cyclically
i.e. $j_t = 1, 2, \dots, d, 1, 2, \dots, d, \dots$

SCD: choose j_t randomly

Block CD: choose a small set of coordinates at each t to update

Step length chosen as usual

(PROJECTED) COORDINATE DESCENT

1. For $t = 0, 1, \dots$

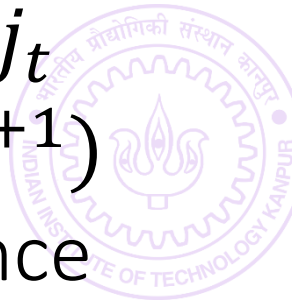
1. Select a coordinate $j_t \in [d]$

2. Let $\mathbf{u}_{j_t}^{t+1} \leftarrow \mathbf{x}_{j_t}^t - \eta_t \cdot \nabla_{j_t} f(\mathbf{x}^t)$

3. Let $\mathbf{u}_j^{t+1} \leftarrow \mathbf{x}_j^t$ for $j \neq j_t$

4. Project $\mathbf{x}^{t+1} \leftarrow \Pi_{\mathcal{C}}(\mathbf{u}^{t+1})$

5. Repeat until convergence



Coordinate Descent

5

Sometimes we are able to optimize completely along a given variable (even if constraints are there) – called *coordinate minimization (CM)*

Will see an example of this today

Coordinate descent popularly applied to dual problems (SDCA – stochastic dual coordinate ascent)

Recall: dual problem is a maximization problem – hence ascent not descent

Beneficial if $n \gg d$ which is quite often the case

CD/CA/CM applied to primal problems too if $d \gg n$ since in these cases, SGD takes $\mathcal{O}(d)$ time per step but SCD usually takes only $\mathcal{O}(n)$

Achieving this speedup can require some clever bookkeeping

SCD naively applied in these cases will take $\mathcal{O}(nd)$ time per step



SDCA for the CSVM Problem

6

Consider the version that does not have a bias

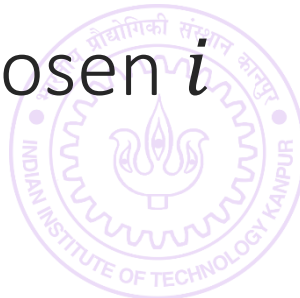
$$\begin{aligned} & \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \boldsymbol{\alpha}_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \boldsymbol{\alpha}_i \boldsymbol{\alpha}_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \right\} \\ & \text{s.t. } \boldsymbol{\alpha}_i \in [0, C] \text{ for all } i \in [n] \end{aligned}$$

*The version with bias has constraint $\sum_{i=1}^n \boldsymbol{\alpha}_i y^i = 0$ that links all $\boldsymbol{\alpha}_i$ together
Cannot update a single $\boldsymbol{\alpha}_i$ without disturbing all the others ☹️*

Need a more involved algorithm Sequential Minimal Optimization (SMO) by John Platt to solve the version with a bias – update two $\boldsymbol{\alpha}_i$ at a time!

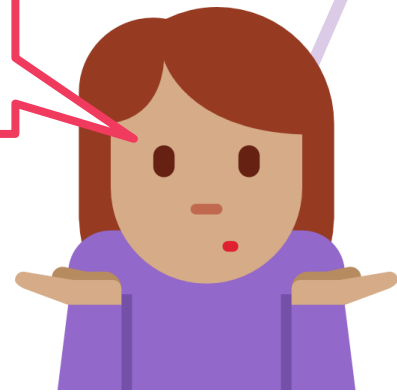
To apply SCA, at each step, we need to choose one $i \in [n]$

In this case, possible to completely maximize objective w.r.t. chosen i



SDCA for the CSVM P

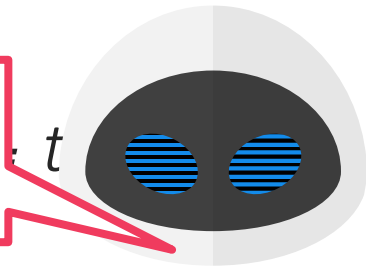
Does this mean I need to choose one data point at each time step?



Consider the version that does not have a bias

$$\begin{aligned} & \underset{\alpha \in \mathbb{R}^n}{\operatorname{argmax}} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \right\} \\ & \text{s.t. } \alpha_i \in [0, C] \text{ for all } i \in [n] \end{aligned}$$

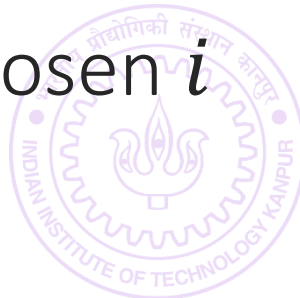
The version without a bias cannot be solved by coordinate ascent. Yes, coordinate ascent in the dual looks a lot like stochastic gradient descent in the primal! Both work with a single data point at a time. Cannot update a single α_i without disturbing all the others ☹



Need a more involved algorithm Sequential Minimal Optimization (SMO) by John Platt to solve the version with a bias – update two α_i at a time!

To apply SCA, at each step, we need to choose one $i \in [n]$

In this case, possible to completely maximize objective w.r.t. chosen i



SDCA for the CSVM Problem

8

$$\begin{aligned} & \operatorname{argmax}_{\boldsymbol{\alpha} \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \boldsymbol{\alpha}_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \boldsymbol{\alpha}_i \boldsymbol{\alpha}_j y^i y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \right\} \\ & \text{s.t. } \boldsymbol{\alpha}_i \in [0, C] \text{ for all } i \in [n] \end{aligned}$$

Concentrating on just the terms that involve $\boldsymbol{\alpha}_i$ we get

$$\begin{aligned} & \operatorname{argmax}_{\boldsymbol{\alpha}_i \in \mathbb{R}} \left\{ \boldsymbol{\alpha}_i - \frac{1}{2} \boldsymbol{\alpha}_i^2 \|\mathbf{x}^i\|_2^2 - \boldsymbol{\alpha}_i y^i \cdot \sum_{j \neq i} \boldsymbol{\alpha}_j y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle \right\} \\ & \text{s.t. } \boldsymbol{\alpha}_i \in [0, C] \end{aligned}$$

Renaming $x = \boldsymbol{\alpha}_i$, $q = \|\mathbf{x}^i\|_2^2$, $p = y^i \cdot \sum_{j \neq i} \boldsymbol{\alpha}_j y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$, we get

$$\operatorname{argmin}_x \frac{1}{2} q x^2 - x(1 - p) \text{ s.t. } x \in [0, C]$$

Solution is very simple: find unrestricted minimum i.e. $\tilde{x} = (1 - p)/q$

If $\tilde{x} \in [0, C]$, solution is \tilde{x} elif $\tilde{x} < 0$, solution is 0, else solution is C



SDCA for the

Warning: in general, finding an unconstrained solution and doing a projection step **does not** give a true solution



$$\operatorname{argmax}_{\alpha \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j v^i v^j \langle x^i, x^j \rangle \right\}$$

$$\text{s.t. } \alpha_i \in [0, C]$$

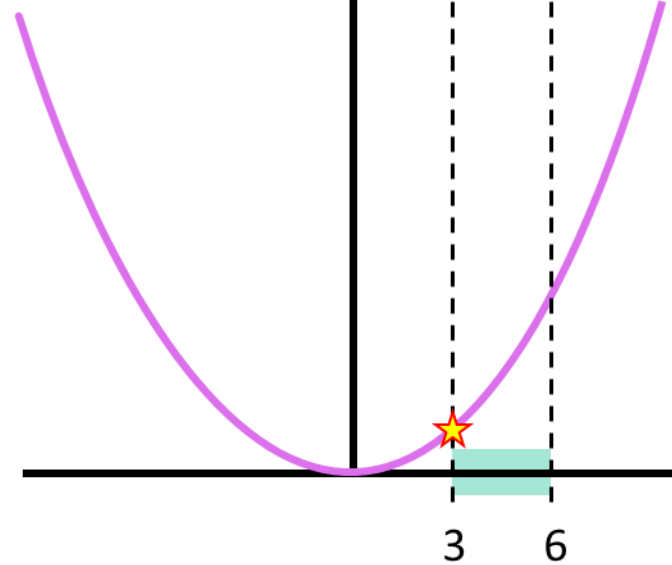
Concentrating on

$$\operatorname{argmax}_{\alpha_i \in \mathbb{R}} \left\{ \alpha_i - \frac{1}{2} \sum_{j=1}^n \alpha_j v^j \langle x^i, x^j \rangle \right\}$$

$$\text{s.t. } \alpha_i \in [0, C]$$

Renaming $x =$

$$\begin{aligned} \min_x & x^2 \\ \text{s.t. } & x \geq 6 \\ & \text{and } x \leq 3 \end{aligned}$$



$\langle x^j, \cdot \rangle$, we get

$$\operatorname{argmin}_x \frac{1}{2} q x^2 - x(1-p) \text{ s.t. } x \in [0, C]$$

Solution is very simple: find

If $\tilde{x} \in [0, C]$, solution is \tilde{x}

Indeed! In this special case, our objective had a nice property called *unimodality* which is why this trick works – it won't work in general

$/q$



Speeding up SDCA computations

10

All that is left is to find how to compute p, q for our chosen i

$q = \|\mathbf{x}^i\|_2^2$ can be easily precomputed for all data points

However, $p = y^i \cdot \sum_{j \neq i} \alpha_j y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle$ needs $\mathcal{O}(nd)$ time to compute ☹

... only if done naively. Recall that we always have $\mathbf{w} = \sum_{i=1}^n \alpha_i y^i \cdot \mathbf{x}^i$ for the CSVM (even if we have bias and slack variables)

Thus, $\sum_{j \neq i} \alpha_j y^j \langle \mathbf{x}^i, \mathbf{x}^j \rangle = \mathbf{w}^\top \mathbf{x}^i - \alpha_i y^i \|\mathbf{x}^i\|_2^2 = \mathbf{w}^\top \mathbf{x}^i - \alpha_i y^i q$

All we need to do is create (and update) the \mathbf{w} vector in addition to the α vector and we would be able to find p in just $\mathcal{O}(d)$ time 😊



Loss Functions

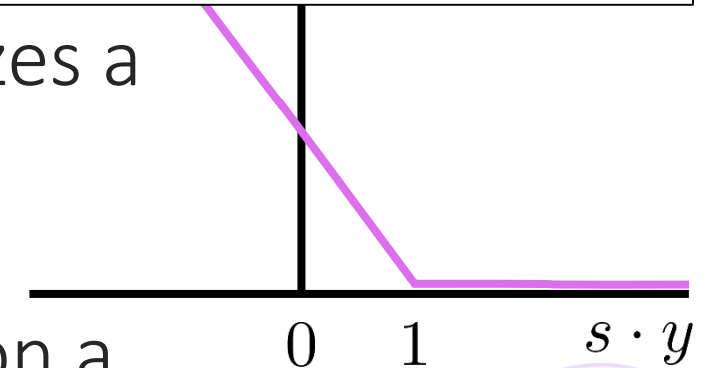
11

Loss functions are widely used in ML to encode what behaviour we desire from our model

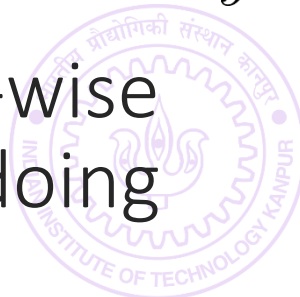
Loss functions penalize undesirable behaviour and when we use an optimizer like SGD or SDCA, we (hopefully) end up with a model which has desirable behaviour

$$\ell_{\text{hinge}}(y^i \cdot \mathbf{w}^\top \mathbf{x}^i) = [1 - y^i \cdot \mathbf{w}^\top \mathbf{x}^i]_+$$

The hinge function is a loss function which penalizes a model if it either misclassifies a data point or else correctly classifies it with insufficient margin



Note: hinge loss tells us how well is model doing on a single data point (\mathbf{x}^i, y^i) . We take sum/avg of these datapoint-wise loss values on the entire training set to see how well is model doing on the entire training set i.e. look at $\sum_{i=1}^n \ell_{\text{hinge}}(y^i \cdot \mathbf{w}^\top \mathbf{x}^i)$



Other Classification Loss Functions

12

Squared Hinge loss Function:

$$\ell_{\text{sq-hinge}}(y^i \cdot \mathbf{w}^\top \mathbf{x}^i) = [1 - y^i \cdot \mathbf{w}^\top \mathbf{x}^i]_+^2$$

Popular since it is differentiable – no kinks

Logistic Loss Function:

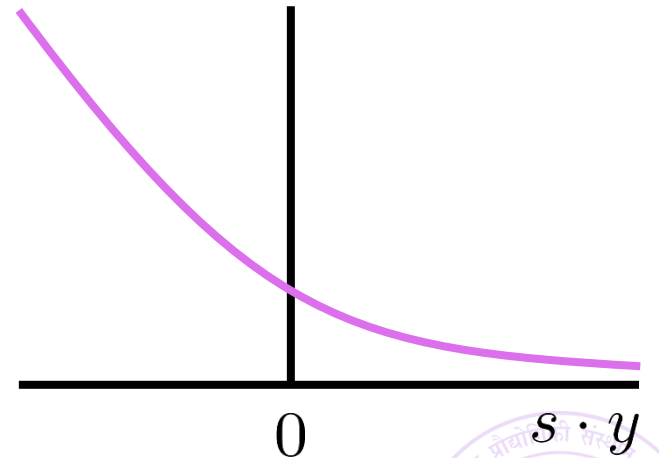
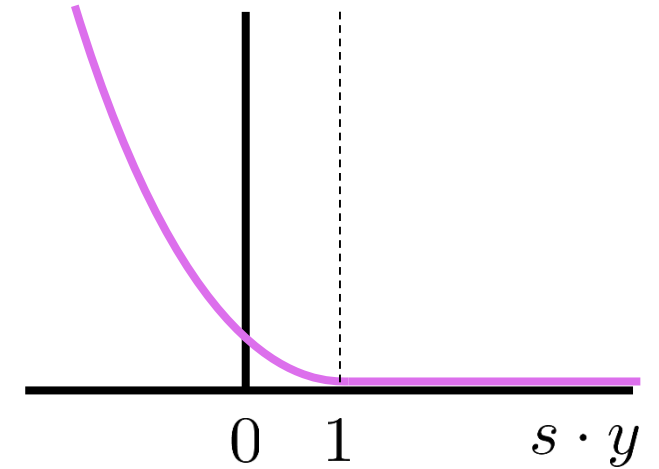
$$\ell_{\text{logistic}}(y^i \cdot \mathbf{w}^\top \mathbf{x}^i) = \ln(1 + \exp(-y^i \cdot \mathbf{w}^\top \mathbf{x}^i))$$

Popular, differentiable

Related to the cross-entropy loss function

Some loss functions e.g. hinge, can be derived in a geometric way, others e.g. logistic, can be derived probabilistically

However, some e.g. squared hinge, are directly proposed by ML experts as they have nice properties – no separate “intuition” for these 😊



Logistic Regression

13

If we obtain a model by solving the following problem

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \cdot \sum_{i=1}^n \ln(1 + \exp(-y^i \cdot \mathbf{w}^\top \mathbf{x}^i))$$

then we would be doing *logistic regression*

Note that we simply replaced hinge loss with logistic loss here

Warning: logistic regression solves a binary classification problem, not a regression problem. The name is misleading

Can be solved in primal by using GD/SGD/MB

Can be solved in dual using SDCA as well

Dual derivation has to be done a bit cleverly here since no constraints 😊

Will revisit logistic regression very soon and derive the loss function



Regression Problems

14

In binary classification, we have to predict one of two classes for each data point i.e. $\mathbf{x} \mapsto \{-1, 1\}$

In regression, we have to predict a real value i.e. $\mathbf{x} \mapsto \mathbb{R}$

Training data looks like $\{(\mathbf{x}^i, y^i)\}_{i=1}^n$ where $\mathbf{x}^i \in \mathbb{R}^d, y \in \mathbb{R}$

For example, predict price of a stock, predict *change* in price of a stock, predict test scores of a student can be solved using regression

Let us look at a few ways to solve regression problem as well as loss functions for regression problems

Recall: logistic regression is not a way to perform regression, it is a way to perform binary classification



Solving Regression Problems via kNN

15

Store all training data as the “model”

For a given test data point \mathbf{x}^t , find its k nearest neighbours

May use Euclidean/learned metric to define neighbours

Suppose neighbours are $(\mathbf{x}^{i_1}, y^{i_1}), \dots, (\mathbf{x}^{i_k}, y^{i_k})$

Predict the average score of neighbours i.e. $\hat{y}^t = \frac{1}{k} \sum_{j=1}^k y^{i_j}$



Solving Regression Problems via DT

16

Need a different notion of “purity”

Hopeless to expect y^i values to repeat at a node and use old notion of purity

May call a node pure if y values in that node are “close” to each other

May define this in many ways: given $(\mathbf{x}^1, y^1), \dots, (\mathbf{x}^n, y^n)$ at a node

Purity definition 1: $\sum_{i=1}^n \sum_{j=1}^n |y^i - y^j|$

Purity definition 2: $\sum_{i=1}^n \sum_{j=1}^n (y^i - y^j)^2$

Purity definition 3: Let $\bar{y} = \frac{1}{n} \sum_{i=1}^n y^i$ and use $\sum_{i=1}^n |y^i - \bar{y}|$ as purity

Purity definition 4: Let $\bar{y} = \frac{1}{n} \sum_{i=1}^n y^i$ and use $\sum_{i=1}^n (y^i - \bar{y})^2$ as purity

Defn 4 is related to *variance*. Defn 3 is related to *absolute deviation*

Possible leaf action: predict average y for training points at that leaf



Loss Functions for Regression Problems

17

Can use linear models to solve regression problems too i.e. learn a (\mathbf{w}, b) and predict score for test data point \mathbf{x}^t as $\hat{y}^t = \mathbf{w}^\top \mathbf{x}^t + b$

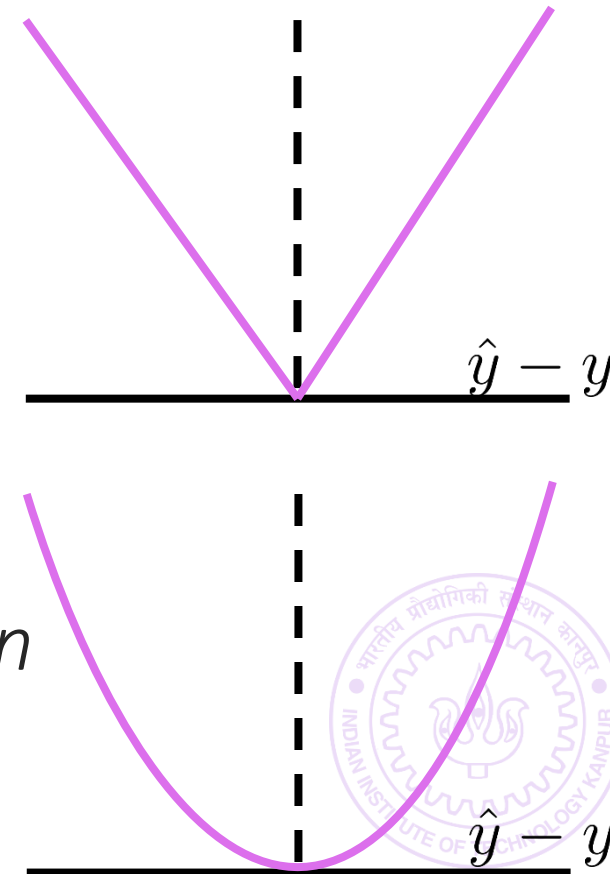
Need loss functions that define what they think is bad behaviour

Absolute Loss: $\ell_{\text{abs}}(\hat{y}, y) = |\hat{y} - y|$

A model is doing badly if \hat{y} is either much larger than y or much smaller than y

Squared Loss: $\ell_{\text{sq}}(\hat{y}, y) = (\hat{y} - y)^2$

A model is doing badly if \hat{y} is either much larger than y or much smaller than y . Also I want the loss function to be differentiable so that I can take gradients etc



Loss Functions for Regression Problems

18

$$\ell_{\epsilon}(y, \hat{y}) = \begin{cases} (y - \hat{y} - \epsilon)^2 & \text{if } \hat{y} < y - \epsilon \\ 0 & \text{if } \hat{y} - y \in [-\epsilon, \epsilon] \\ (y - \hat{y} + \epsilon)^2 & \text{if } \hat{y} > y + \epsilon \end{cases}$$

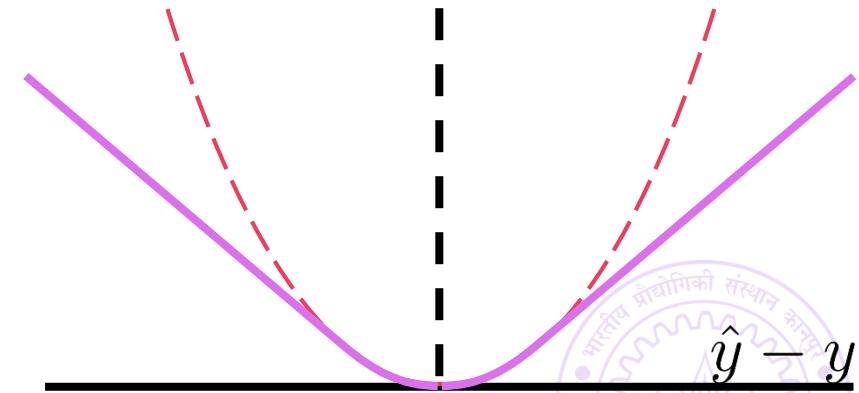
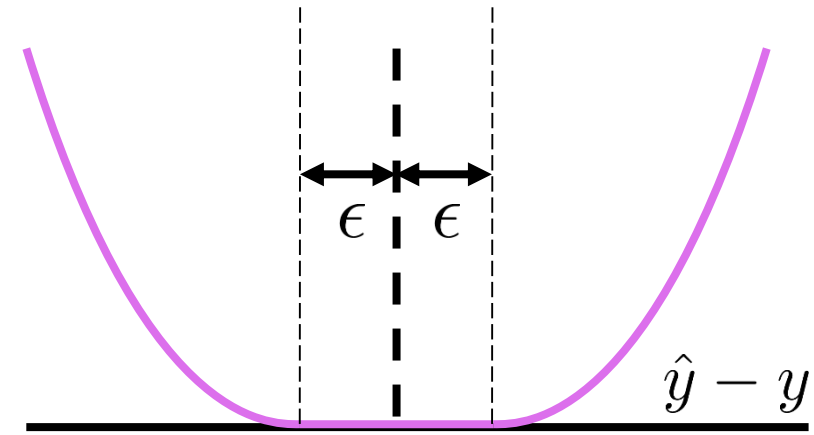
Vapnik's ϵ -insensitive Loss:

If a model is doing slightly badly, don't penalize it at all, else penalize it as squared loss. Ensure a differentiable function

$$\ell_{\delta}(y, \hat{y}) = \begin{cases} (\hat{y} - y)^2 & \text{if } |\hat{y} - y| \leq \delta \\ \delta \cdot |y - \hat{y}| & \text{if } |\hat{y} - y| \geq \delta \end{cases}$$

Huber Loss:

If a model is doing slightly badly, penalize it as squared loss, if doing very badly, penalize it as absolute loss. Ensure a differentiable function



Loss functions have to be chosen according to needs of the problem and experience. Eg. Huber loss popular if some data points are corrupted. However, some loss functions are very popular in ML for various reasons e.g. hinge/cross entropy for binary classification, squared for regression

 $\ell_{\epsilon}(y,$

$$(y - \hat{y} + \epsilon)^2 \quad \text{if } \hat{y} > y + \epsilon$$

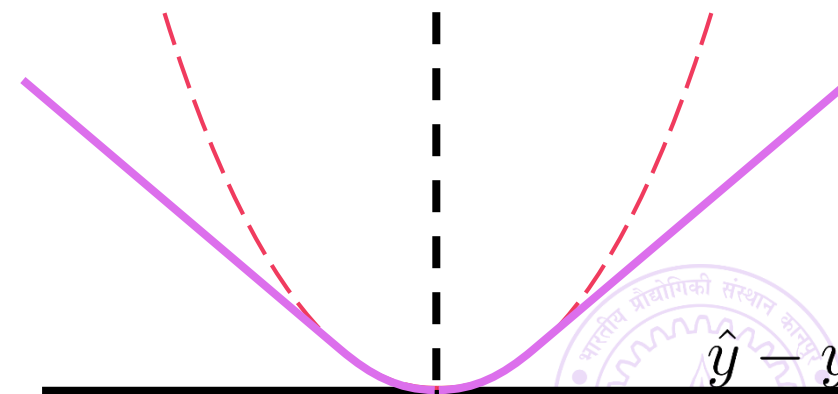
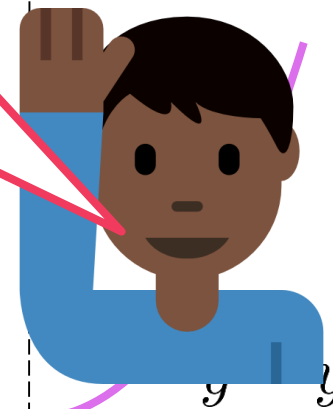
Vapnik's ϵ -insensitive Loss:

If a model is doing slightly badly, don't penalize it at all, else penalize it as squared loss. Ensure a differentiable function

$$\ell_{\delta}(y, \hat{y}) = \begin{cases} (\hat{y} - y)^2 & \text{if } |\hat{y} - y| \leq \delta \\ \delta \cdot |y - \hat{y}| & \text{if } |\hat{y} - y| \geq \delta \end{cases}$$

Huber Loss:

If a model is doing slightly badly, penalize it as squared loss, if doing very badly, penalize it as absolute loss. Ensure a differentiable function



Least Squares Regression

20

Ignore the bias b for sake of notational simplicity

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i)^2$$

Can rewrite as $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \|X\mathbf{w} - \mathbf{y}\|_2^2$ with $X \in \mathbb{R}^{n \times d}, \mathbf{y} \in \mathbb{R}^n$

Can apply first order optimality to obtain a solution (convex objective)

Gradient: $\lambda \cdot \mathbf{w} + 2X^\top (X\mathbf{w} - \mathbf{y})$ (recall: gradient must be $d \times 1$)

Gradient must vanish at minimum so we must have

$$\mathbf{w} = (2X^\top X + \lambda \cdot I)^{-1} (2X^\top \mathbf{y})$$

Takes $\mathcal{O}(d^3)$ time to invert the matrix – may use faster methods (S)GD

Much faster methods available e.g. conjugate gradient method



Learn

Could have used $\arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|_2^2 + C \cdot \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i)^2$ too but it is customary in regression to write the optimization problem a bit differently

Ignore the bias b for sake of notational simplicity

$$\arg \min_{\mathbf{w} \in \mathbb{R}^d} \lambda \|\mathbf{w}\|_2^2 + \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i)^2$$

Can Even here, we may use dual methods like SDCA (liblinear etc do indeed use them) but need to carefully derive the dual since no constraints

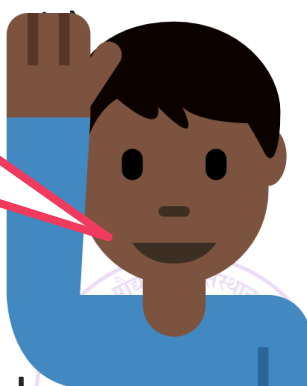
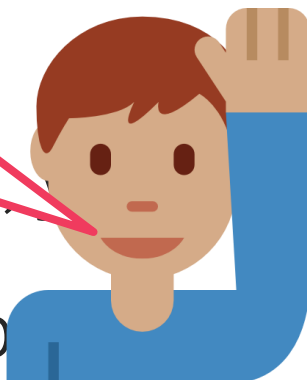
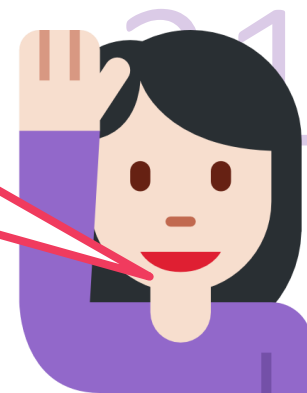
Can apply first order optimality to obtain a solution (convex ob

Grad If we want to use fancier loss functions like Vapnik's loss function, then
Grad cannot apply first order optimality, need to use SGD, SDCA etc methods

$$\mathbf{w} = (2X^\top X + \lambda \cdot I)^{-1} (2X^\top \mathbf{y})$$

Takes $\mathcal{O}(d^3)$ time to invert the matrix – may use faster methods (S)GD

Much faster methods available e.g. conjugate gradient method



Behind the scenes in GD for Least Squares 22

$$f(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i)^2 \text{ (ignore bias } b \text{ for now)}$$

$$\nabla f(\mathbf{w}) = \lambda \cdot \mathbf{w} + 2 \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i) \cdot \mathbf{x}^i$$

$$\mathbf{w}^{\text{new}} = \mathbf{w} - \eta \cdot \nabla f(\mathbf{w}) = (1 - \eta\lambda) \cdot \mathbf{w} - 2\eta \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i) \cdot \mathbf{x}^i$$

Assume $n = 1$ for a moment for sake of understanding

$$\mathbf{w}^{\text{new}} = (1 - \eta\lambda) \cdot \mathbf{w} - 2\eta(\mathbf{w}^\top \mathbf{x}^1 - y^1) \cdot \mathbf{x}^1$$

Small η : $(1 - \eta\lambda)$ is large \Rightarrow do not change \mathbf{w} too much!

If \mathbf{w} does well on (\mathbf{x}^1, y^1) , say $\mathbf{w}^\top \mathbf{x}^1 = y^1$, then $\mathbf{w}^{\text{new}} = (1 - \eta\lambda) \cdot \mathbf{w}$

If \mathbf{w} does badly on (\mathbf{x}^1, y^1) , say $\mathbf{w}^\top \mathbf{x}^1 \gg y^1$, then

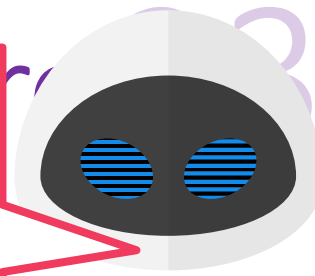
$$\mathbf{w}^{\text{new}} = (1 - \eta\lambda) \cdot \mathbf{w} - 2\eta g \cdot \mathbf{x}^1 \text{ where } g \gg 0$$

$$(\mathbf{w}^{\text{new}})^\top \cdot \mathbf{x}^1 = (1 - \eta\lambda) \cdot \mathbf{w}^\top \mathbf{x}^1 - 2\eta g \cdot \|\mathbf{x}^1\|_2^2$$



Behind the scene

GD and other optimization techniques merely try to obtain models that obey the rules of good behaviour as encoded in the loss function



$$f(\mathbf{w}) = \frac{\lambda}{2} \|\mathbf{w}\|_2^2 + \sum_{i=1}^n \left(\frac{1}{2} (\mathbf{w}^\top \mathbf{x}^i - y^i)^2 \right)$$

$$\nabla f(\mathbf{w}) = \lambda \cdot \mathbf{w} + 2 \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i) \cdot \mathbf{x}^i$$

$$\mathbf{w}^{\text{new}} = \mathbf{w} - \eta \cdot \nabla f(\mathbf{w}) = (1 - \eta\lambda) \cdot \mathbf{w} - 2\eta \sum_{i=1}^n (\mathbf{w}^\top \mathbf{x}^i - y^i) \cdot \mathbf{x}^i$$

Assume $n = 1$ for a moment for sake of understanding

$$\mathbf{w}^{\text{new}} = (1 - \eta\lambda) \cdot \mathbf{w} - 2\eta(\mathbf{w}^\top \mathbf{x}^1 - y^1) \cdot \mathbf{x}^1$$

No change to \mathbf{w} due to the data point (\mathbf{x}^1, y^1)

Small η : $(1 - \eta\lambda)$ is large \Rightarrow do not change \mathbf{w} too much

If \mathbf{w} does well on (\mathbf{x}^1, y^1) , say $\mathbf{w}^\top \mathbf{x}^1 = y^1$, then $\mathbf{w}^{\text{new}} = \mathbf{w}$

If \mathbf{w} does badly on (\mathbf{x}^1, y^1) , say $\mathbf{w}^\top \mathbf{x}^1 \gg y^1$, then

$(\mathbf{w}^{\text{new}})^\top \cdot \mathbf{x}^1$ is smaller than $\mathbf{w}^\top \mathbf{x}^1$ i.e. may be closer to y^1

$$\mathbf{w}^{\text{new}} = (1 - \eta\lambda) \cdot \mathbf{w} - 2\eta g \cdot \mathbf{x}^1 \text{ where } g \gg 0$$

$$(\mathbf{w}^{\text{new}})^\top \cdot \mathbf{x}^1 = (1 - \eta\lambda) \cdot \mathbf{w}^\top \mathbf{x}^1 - 2\eta g \cdot \|\mathbf{x}^1\|_2^2$$

If $\mathbf{w}^\top \mathbf{x}^1 \ll y^1$, GD will try to increase the value