

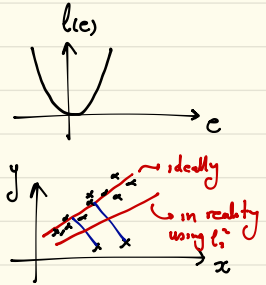
Robust Regression

One point about using l_2^2 cost function is its high intolerance for large errors

$$(y - h_\theta(x))^2 = l_2(y - h_\theta(x)) = e^2$$

What if the dataset contains "outliers" ?

l_2^2 pretends to have "many smaller errors"
than "a few large errors" !



Solution use a more robust cost function to deal with outliers

Let $l(e) = |e|$ (absolute value or l_1 -norm)

$$\min_{\theta} \sum_{i=1}^N |y_i - h_\theta(x_i)| = \sum_{i=1}^N |y_i - \theta^T x_i| = \|Y - X\theta\|_1$$

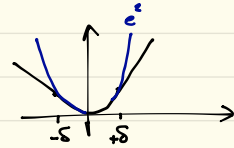
Still convex!

challenge $|e|$ is non-diff at its minimum $\frac{\partial |e|}{\partial e} = \begin{cases} +1 & , e > 0 \\ -1 & , e < 0 \\ \text{undefined} & , e = 0 \end{cases}$

unlike l_2^2 not easy to optimize (no closed-form, no grad desc), Yet convex!

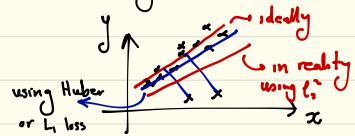
Another objective function that has the advantages of both norms robustness of $|e|$ and differentiability of e^2 is called "Huber loss"

$$l_{\frac{\delta}{2}}(e) \triangleq \begin{cases} e^2/2, & |e| \leq \delta \\ \delta|e| - \frac{\delta^2}{2}, & |e| > \delta \end{cases}$$

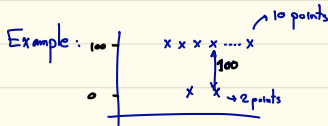


Convex, $\in C^1$ (has a continuous grad) is differentiable everywhere!

⇒ Does it have a closed-form solution?



To optimize, we can use batch/stochastic gradient descent!



Overfitting

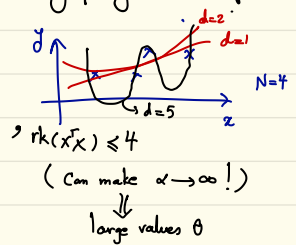
A major problem of learning systems is "overfitting" performing well on training samples but performing poorly on test samples!

$$d=2 \Rightarrow h(x) = \theta_0 x + b \Rightarrow \begin{matrix} \bar{X}^T \theta = \bar{X}^T Y \\ \text{2x1} \quad \text{2x1} \quad \text{2x1} \end{matrix} \xrightarrow{N \geq d} \theta \checkmark$$

$$d=3 \Rightarrow h(x) = \theta_0 x + \theta_1 x^2 + b \Rightarrow \begin{matrix} \bar{X}^T \theta = \bar{X}^T Y \\ \text{3x1} \quad \text{3x1} \quad \text{3x1} \end{matrix} \xrightarrow{N \geq d} \theta \checkmark$$

$$d=5 \Rightarrow h(x) = \theta_0 x + \theta_1 x^2 + \theta_2 x^3 + \theta_3 x^4 + b \xrightarrow{N < d} \bar{X}^T \theta = \bar{X}^T Y$$

$$\hookrightarrow \bar{X}^T v = 0 \quad 3 \times 1 \rightarrow \bar{X}^T (\theta^* + \alpha v) = \bar{X}^T \theta^* = \bar{X}^T Y$$



Treatment: Regularization

$$\min_{\theta} \sum_{i=1}^N \ell(y_i - \theta^T x_i) + \underbrace{\lambda \frac{1}{2} \|\theta\|_2^2}_{\text{regularization term}} \quad \text{reg parameter } > 0$$

$$\ell(\theta) = \|\theta\|_2^2, \quad \ell(c) = c^2 \rightarrow \text{Ridge regression}$$

$$\min_{\theta} \|\gamma - X\theta\|_2^2 + \underbrace{\lambda \|\theta\|_2^2}_{\text{weight decay}} \triangleq J(\theta)$$

$J(\cdot)$ convex & diff wrt θ

$$\frac{\partial J}{\partial \theta} = 2X^T(X\theta - \gamma) + 2\lambda\theta$$

$$(1) \nabla J(\theta) \big|_{\theta} = 0 \rightarrow (\bar{X}^T \bar{X} + \lambda I) \theta = \bar{X}^T Y \rightarrow \theta = (\bar{X}^T \bar{X} + \lambda I)^{-1} \bar{X}^T Y$$

Notice that $\lambda I + \bar{X}^T \bar{X}$ always invertible for $\lambda > 0$.

\rightarrow it is always invertible even when $N < d$

(2) Batch / stochastic GD

\swarrow parameter : learned during training : e.g., θ
 \swarrow hyperparameter : not learned during training, set before, e.g., λ [if we want to learn λ to minimize training then we always pick $\lambda=0$]

λ hyperparameter not learned by the learning alg, is set in advance
 (however, in a full Bayesian method, has a prob for it)

Effect of λ $\lambda \rightarrow 0$ same as LS, overfitting if $d > N$
 $\lambda \rightarrow \infty$ $\hat{\theta} = 0$ ($(\tilde{X}^T \tilde{X} + \lambda I)^{-1} \tilde{X}^T \tilde{y} \approx 0$) \rightarrow No overfit but high training & test error

Question How to set λ ?

Enough to give enough expressive power to model to obtain small training and test error (even $\lambda=10^{-4}$ is better than nothing)

+ Naive way - For each λ , solve $\min_{\theta} \sum_{i=1}^N \ell(y_i - \theta^T x_i) + \lambda f(\theta) \rightarrow \hat{\theta}_{\lambda}$

- compute $\mathcal{E}_{\text{train}}(\lambda) = \sum_{i=1}^N \ell(y_i - \hat{\theta}_{\lambda}^T x_i)$

- take $\lambda^* = \arg \min_{\lambda} \mathcal{E}_{\text{train}}(\lambda)$

$\Rightarrow \lambda^* = 0$ overfitting !!



* Right approach "Hold out" set set aside part of the training data as hold out set

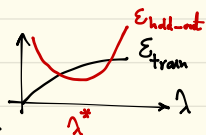


+ For each λ , compute $\hat{\theta}_\lambda$ using D^{tr}

+ " " " , compute $E_{\text{hold-out}}(\lambda) = \sum_{(x_i, y_i) \in D^{ho}} \ell(y_i, -\hat{\theta}_\lambda^T x_i)$

\Rightarrow This can be a good measure of generalization error, i.e., error on unseen data

+ take $\lambda^* = \underset{\lambda}{\operatorname{argmin}} E_{\text{hold-out}}(\lambda)$

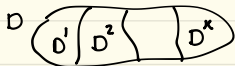


+ run regularized regression on test data using λ^*

\rightarrow This is how we can effectively deal with generalization and overfitting

Challenge if training dataset is small, hold-out set makes training set smaller \rightarrow not good $\hat{\theta}_\lambda$ due to more overfitting!

Solution K-fold cross validation divide training set into K partitions, run training on all but one partition, evaluate on the remaining part.



- For $\ell=1, \dots, K$

- find $\hat{\theta}_\lambda$ using $\bigcup_{j \neq \ell} D^j$, for several values of λ

- compute $E_{\text{hold-out}}^{(\ell)}(\lambda)$ using $D^{(\ell)}$

- compute the average hold-out error as

$$E_{\text{hold-out}}(\lambda) = \frac{1}{K} \sum_{j=1}^K E_{\text{hold-out}}^{(j)}(\lambda)$$

- take $\lambda^* = \underset{\lambda}{\operatorname{argmin}} E_{\text{hold-out}}(\lambda)$

- test on test samples using λ^*

* If $K=N$ (training on all but one data at a time), this is called One hold out CV \rightarrow Expensive, but effective when N small