

1 Linear Models

1.1 Overfitting

- fitting: der Prozess die Parameter einer Modelfunktion y so anzupassen das sie der der Beispieldaten D am besten passen
- overfiiting: "Fitting the data more than is warranted"
- alias besser passen als berechtigt?
- Gründe:
 - zu komplizierte Modelfunktion (zu viele Features)
 - zu wenig Daten in D
 - zu viel Datenrauschen
 - D ist zu biased alias nicht repräsentativ
- Folgen:
 - kleiner Error auf D_{tr} aber großer Error auf D_{test} und IRL
 - loss of inductive Bias
 - increase of variance as a result of sensitivity to noise
- Overfitting finden:
 - Visuell untersuchen für Fälle mit Dimensionen < 3 sonst embedding oder projizieren in kleinere Dimensionen
 - Validieren: wenn $Err_{fit} = Err_{val}(y) - Err_{tr}(y)$ zu groß ist
- Overfitting vermeiden:
 - Early stopping through model selection: nach m schritten überprüfen ob sich Err_{fit} noch verkleinert und stoppen wenn er sich vergrößert
 - Qualität (schlechte Beispiele raus) und / oder Quantität (mehr Daten gleichen Rauschen aus) von D verbessern
 - Manually enforcing a higher bias by using a less complex hypothesis space
 - Regularization (WUHU!)

1.1.1 Well- and Ill-posed problems

A mathematical problem is called well-posed if

1. a solution exists,
2. the solution is unique,
3. the solution's behavior changes continuously with the initial conditions.

Otherwise, the problem is called ill-posed.

1.2 Regularization

Automatic adjustment of the loss function to penalize model complexity. Let $L(\mathbf{w})$ denote a loss function used to optimize the parameters \mathbf{w} of a model function $y(\mathbf{x})$. Regularization introduces a trade-off between model complexity and inductive bias:

$$\mathcal{L}(\mathbf{w}) = L(\mathbf{w}) + \lambda * R(\mathbf{w})$$

where $\lambda \geq 0$ controls the impact of the regularization term $R(\mathbf{w}) \geq 0$. \mathcal{L} is called “objective function”.

1.2.1 Regularized Linear Regression

$$\mathcal{L}(\mathbf{w}) = \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \cdot \vec{\mathbf{w}}^T \vec{\mathbf{w}}$$

Estimate \mathbf{w} by minimizing the residual sum of squares:

$$\hat{\mathbf{w}} = \underset{\mathbf{w} \in \mathbf{R}^{p+1}}{\operatorname{argmin}} \mathcal{L}(\mathbf{w})$$

$$\rightsquigarrow RSS(\mathbf{w}) = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \lambda \mathbf{w}^T \mathbf{w}$$

Ableitung bilden um $RSS(\mathbf{w})$ zu minimieren und man kommt auf:

$$\mathbf{w} = (X^T X + \operatorname{diag}(0, \lambda, \dots, \lambda))^{-1} X^T \mathbf{y}$$

$$\operatorname{diag}(0, \lambda, \dots, \lambda) = \begin{bmatrix} 0 & 0 & 0 & \dots & 0 \\ 0 & \lambda & 0 & \dots & 0 \\ 0 & 0 & \lambda & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \lambda \end{bmatrix}$$

$$\hat{y}(\mathbf{x}_i) = \hat{\mathbf{w}}^T \mathbf{x}_i$$

Um so höher λ um so einfacher ist die Funktion - Regularization archived!

2 Neural Networks

2.1 Perception Learning

Idee: Lass mal ein Gehirn programmieren!

$$y(\mathbf{x}) = 1 \Leftrightarrow \left(\sum_{j=0}^p w_j x_j \right) \geq 0$$

sonst ist $y(\mathbf{x}) = 0$

- wenn $w_0 = -\theta$ und $x_0 = 1$ (canonical form)
- sonst $y(\mathbf{x}) = 1 \Leftrightarrow \left(\sum_{j=1}^p w_j x_j - \theta \right) \geq 0$

2.1.1 PT Algorithm

PT(D, η)

1. *initialize_random_weights(w), t = 0*
2. **REPEAT**
3. $t = t + 1$
4. $(\mathbf{x}, c(\mathbf{x})) = \text{random_select}(D)$
5. $\delta = c(\mathbf{x}) - y(\mathbf{x}) \quad // \quad y(\mathbf{x}) \stackrel{(*)}{=} \text{heaviside}(\mathbf{w}^T \mathbf{x}) \in \{0, 1\}, \quad c(\mathbf{x}) \in \{0, 1\} \rightsquigarrow$
6. $\Delta \mathbf{w} \stackrel{(*)}{=} \eta \cdot \delta \cdot \mathbf{x}$
7. $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$
8. **UNTIL**(*convergence*($c(D), y(D)$) **OR** $t > t_{\max}$)
9. *return(w)*

- If a separating hyperplane between X_0 and X_1 exists, the PT algorithm will converge. If no such hyperplane exists, convergence cannot be guaranteed.
- A separating hyperplane can be found in polynomial time with linear programming. The PT Algorithm, however, may require an exponential number of iterations.
- Classification problems with noise are problematic

2.2 Gradient Descent

- Finde den kürzesten Weg in ein Min/Max über partielle Ableitungen
- The gradient of a function is the direction of steepest ascent or descent.

2.2.1 Linear Regression + Squared Loss

$$L_2(\mathbf{w}) = \frac{1}{2} \cdot \sum_{(\mathbf{x}, c(\mathbf{x})) \in D} (c(\mathbf{x}) - y(\mathbf{x}))^2$$

Jetzt müssen wir für jedes w_i aus \mathbf{w} eine partielle Ableitung machen um den weight vector zu updaten ($\mathbf{w} = \mathbf{w} + \Delta\mathbf{w}$)

$$\Delta\mathbf{w} = \frac{\delta}{\delta w_i} L_2(\mathbf{w}) = \eta \cdot \sum_D (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \mathbf{x}$$

η = learning rate - a small positiv constant - legen wir selbst fest

2.2.2 The Batch Gradient Descent (BGD) Algorithm

BGD(D, η)

1. *initialize_random_weights*(\mathbf{w}), $t = 0$
2. **REPEAT**
3. $t = t + 1$
4. $\Delta\mathbf{w} = 0$
5. **FOREACH** $(\mathbf{x}, c(\mathbf{x})) \in D$ **DO**
6. $\delta = c(\mathbf{x}) - y(\mathbf{x})$ // $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$, $\delta \in \mathbb{R}$.
7. $\Delta\mathbf{w} \stackrel{(*)}{=} \Delta\mathbf{w} + \eta \cdot \delta \cdot \mathbf{x}$ // $\delta \cdot \mathbf{x}$ is the derivative of ℓ
8. **ENDDO**
9. $\mathbf{w} = \mathbf{w} + \Delta\mathbf{w}$ // $\Delta\mathbf{w}$ is $-\eta \cdot \nabla L_2(\mathbf{w})$ here.
10. **UNTIL**(*convergence*($c(D), y(D)$) **OR** $t > t_{\max}$)
11. *return*(\mathbf{w})

- wichtig: immer wenn irgendwo $\mathbf{w}^T \mathbf{x}$ steht haben wir $x_0 = 1$ zu \mathbf{x} hinzugefügt

- funktionsweise BGD. wir berechnen über die Ableitung in welche Richtung wir müssen und gehen dann einen Schritt der große η
- die "convergence" schaut ob der Squared Loss noch größer als ein ε ist (das wir auch festlegen)
- BGD ist nicht der schnellste (aber sehr einfach) (Newton-Raphson algorithm, BFGS algorithm sind z.B. schneller)
- BGD nimmt den global loss: loss of all examples in D ("batch gradient descent") (Schritt in Richtung die für alle Punkte am besten ist)
- man kann auch den (squared) loss in Bezug auf einzelne Beispiele nehmen (pointwise loss) (dann gehts halt im Zickzack runter)
berechnet sich dann $\ell_2(c(\mathbf{x}), y(\mathbf{x})) = \frac{1}{2}(c(\mathbf{x}) - \mathbf{w}^T \mathbf{x})^2$
- bzw. die weight adaptation: $\Delta \mathbf{w} = \eta \cdot (c(\mathbf{x}) - \mathbf{w}^T \mathbf{x}) \cdot \mathbf{x}$
- für BGD_σ wird Zeile 9 zu
 $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w} + \eta \cdot 2\lambda \cdot \left(\frac{0}{\mathbf{w}}\right)$

2.2.3 The Incremental Gradient Descent IGD Algorithm

IGD(D, η)

1. *initialize_random_weights(w), t = 0*
2. **REPEAT**
3. $t = t + 1$
4. **FOREACH** ($\mathbf{x}, c(\mathbf{x}) \in D$) **DO**
5. $\delta = c(\mathbf{x}) - y(\mathbf{x})$ // $y(\mathbf{x}) \stackrel{(*)}{=} \mathbf{w}^T \mathbf{x}$, $\delta \in \mathbb{R}$.
6. $\Delta \mathbf{w} \stackrel{(*)}{=} \eta \cdot \delta \cdot \mathbf{x}$ // $\delta \cdot \mathbf{x}$ is the derivative of $\ell_2(c(\mathbf{x}), y(\mathbf{x}))$
7. $\mathbf{w} = \mathbf{w} + \Delta \mathbf{w}$
8. **ENDDO**
9. **UNTIL**(*convergence*($c(D), y(D)$) **OR** $t > t_{\max}$)
10. *return(w)*

- kleinere Schritte als BGD
- can better avoid getting stuck in a local minimum of the loss function than BGD

2.2.4 Logistic Regression + Logistic Loss + Regularization

Wie oben nur mit neuer Formel für $\Delta\mathbf{w}$:

$$\Delta\mathbf{w} = \eta \cdot \sum_D (c(\mathbf{x}) - \sigma(\mathbf{w}^T \mathbf{x})) \cdot \mathbf{x} - \eta \cdot 2\lambda \cdot \left(\frac{0}{\mathbf{w}} \right)$$