

Regularizacija dubokih modela

Josip Krapac i Siniša Šegvić

- Regularizacija
- Regularizacija normom vektora parametara modela
- Regularizacija generiranjem podataka i unošenjem šuma
- Regularizacija ranim zaustavljanjem
- Polunadzirano i višezadačno učenje, dijeljenje parametara
- Regularizacija baggingom i dropout

Regularizacija: pregled

Glavni izazov u strojnom učenju: osigurati da model radi dobro ne samo na podacima za učenje nego i na novim podacima.

Tehnike regularizacije: smanjenje greške na skupu za testiranje, uz moguće povećanje greške na skupu za učenje.

Duboki modeli omogućavaju primjenu raznih tehnika regularizacije.

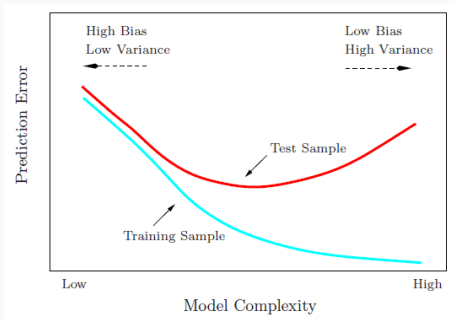
Jedan od najvažnijih otvorenih izazova: vrlo aktivno područje istraživanja.

Regularizacija: pregled

Tehnike regularizacije u strojnom učenju povećavaju pristranosti i smanjuju varijancu modela.

Krajnji cilj: otežavanje **prenaučenosti**.

U praksi: najbolju generalizacijsku pogrešku postiže **model vrlo velikog kapaciteta** na koji su primjenjene **odgovarajuće tehnike regularizacije**.



- Regularizacija
- Regularizacija normom vektora parametara modela
- Regularizacija generiranjem podataka i unošenjem šuma
- Regularizacija ranim zaustavljanjem
- Polunadzirano i višezadačno učenje, dijeljenje parametara
- Regularizacija baggingom i dropout

Regularizacija normom vektora parametara modela

Jedna od najstarijih metoda regularizacije: modificirati gubitka $J(\Theta; \mathbf{X}, \mathbf{y})$ dodavanjem norme vektora parametara $\Omega(\Theta)$:

$$\tilde{J}(\Theta; \mathbf{X}, \mathbf{y}) = J(\Theta; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\Theta)$$

$\alpha \in [0, \infty]$ određuje relativni doprinos regularizatora Ω .

Minimizacija regularizirane funkcije gubitka \tilde{J} smanjuje i J i Ω .

Obično regulariziramo samo **težine**, tj. Ω ne djeluje na **pomak**.

Odabirom funkcije Ω **preferiramo** određene klase modela.

Intuicija: regularizator povlači vektor težina \mathbf{w} prema ishodištu.

Regularizacija L_2 normom vektora parametara modela

Promotrimo funkciju cilja regulariziranu s $\Omega(\Theta) = \frac{1}{2} \|\mathbf{w}\|_2^2$

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = J(\mathbf{w}; \mathbf{X}, \mathbf{y}) + \frac{\alpha}{2} \|\mathbf{w}\|_2^2$$

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}) + \alpha \mathbf{w}$$

Korak gradijentnog spusta sada je ("weight decay"):

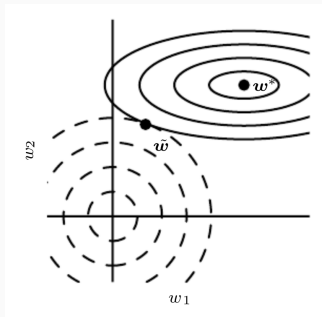
$$\mathbf{w}^{t+1} = \mathbf{w}^t - \epsilon \alpha \mathbf{w}^t - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})|_{\mathbf{w}=\mathbf{w}^t}$$

$$\mathbf{w}^{t+1} = (1 - \epsilon \alpha) \mathbf{w}^t - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})|_{\mathbf{w}=\mathbf{w}^t}$$

Regularizacija L_2 normom vektora parametara modela

Razvijmo ne-regulariziranu funkciju gubitka J u Taylorov red oko minimuma $\mathbf{w}^* = \arg \min_{\mathbf{w}} J(\mathbf{w})$:

$$\hat{J}(\mathbf{w}) = J(\mathbf{w}^*) + \frac{1}{2}(\mathbf{w} - \mathbf{w}^*)^\top \mathbf{H}(\mathbf{w} - \mathbf{w}^*)$$



Regularizacija L_2 normom vektora parametara modela

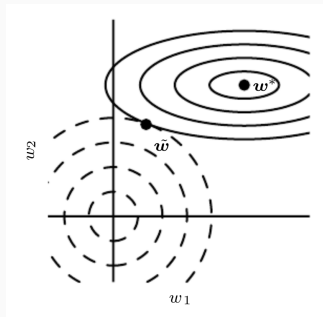
Dodajmo regularizacijski član i pogledajmo gdje se pomakne minimum

$$\nabla_w \left(\hat{J}(\tilde{\mathbf{w}}) + \frac{\alpha}{2} \|\tilde{\mathbf{w}}\|_2^2 \right) = 0$$

$$\nabla_w \hat{J}(\tilde{\mathbf{w}}) + \alpha \tilde{\mathbf{w}} = 0$$

$$\mathbf{H}(\tilde{\mathbf{w}} - \mathbf{w}^*) + \alpha \tilde{\mathbf{w}} = 0$$

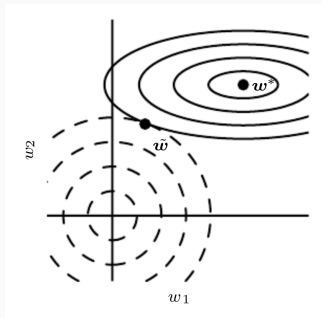
$$\tilde{\mathbf{w}} = (\mathbf{H} + \alpha \mathbf{I})^{-1} \mathbf{H} \mathbf{w}^*$$



Regularizacija L_2 normom vektora parametara modela

Pogledajmo što se dešava kada α raste. Uvid je lakši u prostoru razapetom svojstvenim vektorima (Q) matrice $H = Q\Lambda Q^\top$

$$\begin{aligned}\tilde{w} &= (Q\Lambda Q + \alpha I)^{-1} Q\Lambda Q w^* \\ &= \left[Q(\Lambda + \alpha I) Q^\top \right]^{-1} Q\Lambda Q^\top w^* \\ &= Q(\Lambda + \alpha I)^{-1} \Lambda Q^\top w^* \\ &= Q \cdot \text{diag}\left(\frac{\lambda_i}{\lambda_i + \alpha}\right) \cdot Q^\top w^*\end{aligned}$$



Regularizacija L_2 normom vektora parametara modela

Pogledajmo doprinos projekcije \mathbf{w}^* na i -ti svojstveni vektor \mathbf{H} :

$$\tilde{\mathbf{w}}^{(i)} = \mathbf{q}_i \frac{\lambda_i}{\lambda_i + \alpha} \mathbf{q}_i^\top \mathbf{w}^*$$

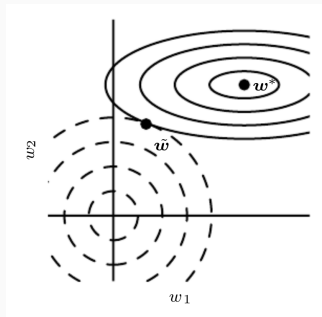
Vidimo da $\tilde{\mathbf{w}}^{(i)}$ postaje tim različitiji od \mathbf{w}_i^* što je λ_i manji, odnosno što je funkcija cilja manje strma u smjeru \mathbf{q}_i

Regularizacija L_2 normom vektora parametara modela

Za slučaj dijagonalne matrice \mathbf{H} ($\mathbf{Q} = \mathbf{I}$):

$$\tilde{\mathbf{w}} = \text{diag}\left(\frac{\lambda_i}{\lambda_i + \alpha}\right) \cdot \mathbf{w}^*$$

$$\tilde{w}_i = \frac{\lambda_i}{\lambda_i + \alpha} w_i^*$$



Vidimo da \tilde{w}_i postaje tim različitiji od w_i što je λ_i manji, odnosno što se \hat{J} manje mijenja po odgovarajućoj osi koordinatnog sustava

Regularizacija identificira parametre koji ne utječu na funkciju cilja i priteže ih prema ishodištu

Regularizacija L_2 normom vektora parametara modela

Pogledajmo utjecaj L_2 regularizacije na linearnu regresiju sa srednjom kvadratnom pogreškom odstupanja kao funkcijom cilja.

Rješenje u ne-regulariziranom slučaju: $\mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$

Rješenje u regulariziranom slučaju: $\mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y}$

Matrica $\mathbf{X}^\top \mathbf{X}$ je proporcionalna kovarijacijskoj matrici $\frac{1}{n} \mathbf{X}^\top \mathbf{X}$.

Efekt regularizacije: prividno povećavanje varijance podataka. Kao da smo oko svakog podatka \mathbf{x} generirali nove podatke izvlačenjem iz normalne distribucije s srednjom vrijednosti koja odgovara podatku \mathbf{x} i varijancom $\alpha \mathbf{I}$.

Regularizacija L_1 normom vektora parametara modela

L_1 regularizator: $\Omega(\Theta) = \sum_i |w_i| = \|\mathbf{w}\|_1$.

Promotrimo gradijent regularizirane funkcije cilja:

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}) + \alpha \text{sgn}(\mathbf{w})$$

Korak gradijentnog spusta je:

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \epsilon \alpha \text{sgn}(\mathbf{w}^t) - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})|_{\mathbf{w}=\mathbf{w}^t}$$

Doprinos regularizacije ovisi samo o predznaku \mathbf{w} .

Regularizacija L_1 normom vektora parametara modela

Isto kao i za L_2 normu: promatramo rastav u Taylorov red u minimumu ne-regularizirane funkcije gubitka, dodamo regularizacijski član i promatramo gdje se pomakne minimum.

Dodatno, pretpostavljamo da je matrica $\mathbf{H} = \mathbf{\Lambda}$ dijagonalna:

$$\tilde{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_i \left(\frac{1}{2} \lambda_i (w_i - w_i^*)^2 + \alpha |w_i| \right)$$

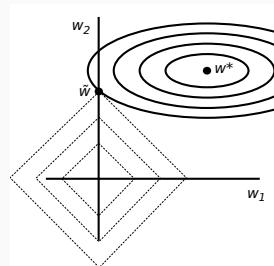
Za ovaj optimizacijski problem postoji analitičko rješenje:

$$\tilde{w}_i = \operatorname{sgn}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{\lambda_i}, 0 \right\}$$

Regularizacija L_1 normom vektora parametara modela

Rješenje L1-regulariziranog kvadratnog problema:

$$\tilde{w}_i = \text{sgn}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{\lambda_i}, 0 \right\}$$

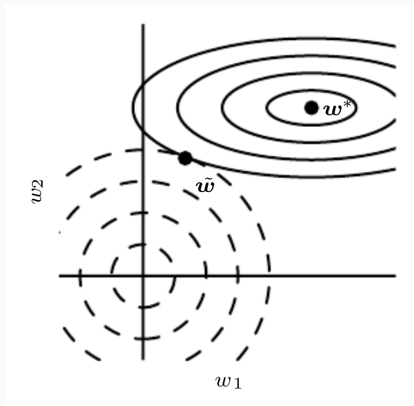


L_1 regularizacija vodi na **rijetke modele**: modeli za koje su neke vrijednosti parametara 0.

L_1 regularizacija istovremeno uči model i obavlja selekciju/eliminaciju značajki.

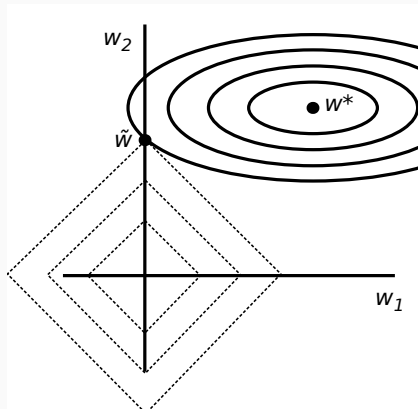
Usporedba L_1 i L_2 regularizacije

L_2



$$\tilde{w}_i = \frac{\lambda_i}{\lambda_i + \alpha} w_i^*$$

L_1



$$\tilde{w}_i = \text{sgn}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{\lambda_i}, 0 \right\}$$

Obje regularizacije guraju w_i prema ishodištu, i to tim jače što je w_i manje važan za podatkovni gubitak (tj. tim jače što je λ_i - manji)

Norme vektora parametara kao ograničenja optimizacijskog postupka

Alternativni pogled: regularizacija osigurava da parametri budu unutar kugle $\Omega(\Theta) < k$.

Radius k ovisan je o parametru α : veći α znači manji k i obratno.

Oblik kugle ovisi o korištenoj normi (npr. L_1 , L_2)

Regularizaciju možemo izraziti i direktno preko k . Modifikacija optimizacijskog postupka:

- napravimo korak gradijentnog spusta i
- projiciramo ažurirani Θ na najbližu točku koja zadovoljava ograničenje $\Omega(\Theta) < k$.

Norme vektora parametara kao ograničenja optimizacijskog postupka

Ova formulacija ne mijenja funkciju gubitka.

- Promjena funkcije gubitka može uzrokovati da optimizacijski postupak zapne u dijelu prostora koji odgovara malim vrijednostima parametara Θ .
- U ovakvoj formulaciji to se izbjegava budući se projekcija obavlja tek kada vektor parametara naraste dovoljno da izađe iz kugle.

Ova formulacija omogućava stabilniji optimizacijski postupak:

- kod korištenja velikih koraka učenja moguće je da postupak počne divergirati (vektor parametara počinje rasti) zbog pozitivne povratne veze.
- U ovakvoj formulaciji nekontroliran rast nije moguć.

- Regularizacija
- Regularizacija norme vektora parametara modela
- Regularizacija generiranjem podataka i unošenjem šuma
- Regularizacija ranim zaustavljanjem
- Polunadzirano i višezadačno učenje, dijeljenje parametara
- Regularizacija baggingom i dropout

Regularizacija generiranjem podataka

Podatci za učenje mogu se promatrati kao odličan regularizator.

Neki problemi dozvoljavaju jednostavno generiranje podataka za učenje modifikacijom postojećih podataka za učenje.

Ova tehnika se pokazala jako uspješnom kod raspoznavanja slika (translatiranje, skaliranje, rotacija) i za raspoznavanje govora.

Pretpostavka: modifikacija ne utječe na ishod predikcije.

Danas je modificiranje podataka popularno za učenje reprezentacija sa zamjenskim gubitkom na neoznačenim podacima [kolesnikov19cvpr]

- npr. model se uči da pogodi parametar modifikacije (rotaciju)...
- ...ili da razlikuje modificirani podatak od drugih podataka itd.

Regularizacija unošenjem šuma

Unošenje šuma u model također ima regularizacijski efekt.

Šum možemo primijeniti na:

- ulaz modela (podatke iz skupa za učenje),
- reprezentacije u skrivenim slojevima,
- parametre modela,
- oznake u skupu za učenje (eng. label smoothing).
 - Ako pretpostavimo da je vjerojatnost točnog označavanja $1 - \epsilon$ onda one-hot oznak transformiramo:

$$[0, 0, \dots, 1, 0, \dots, 0] \rightarrow \left[\frac{\epsilon}{k}, \frac{\epsilon}{k}, \dots, 1 - \frac{k-1}{k}\epsilon, \frac{\epsilon}{k}, \dots, \frac{\epsilon}{k} \right]$$

- Regularizacija
- Regularizacija norme vektora parametara modela
- Regularizacija generiranjem podataka i unošenjem šuma
- Regularizacija ranim zaustavljanjem
- Polunadzirano i višezadačno učenje, dijeljenje parametara
- Regularizacija baggingom i dropout

Regularizacija ranim zaustavljanjem

Broj epoha je hiper-parametar koji možemo efikasno, inkrementalno validirati u jednoj epizodi učenja

Jedino što moramo promijeniti u odnosu na osnovnu proceduru:

- povremeno evaluirati model na skupu za validaciju
- čuvati parametre koji postižu najbolju validacijsku točnost
- usporavanje učenja uslijed evaluiranja možemo ublažiti:
 - evaluacijom na drugom procesoru
 - smanjenjem skupa za evaluaciju
 - rjeđom validacijom

Nedostatak: moramo imati skup za validaciju, što znači da prije konačnog testiranja trebamo ponoviti učenje na trainval skupu

Regularizacija ranim zaustavljanjem

Algorithm 7.1 The early stopping meta-algorithm for determining the best amount of time to train. This meta-algorithm is a general strategy that works well with a variety of training algorithms and ways of quantifying error on the validation set.

Let n be the number of steps between evaluations.

Let p be the “patience,” the number of times to observe worsening validation set error before giving up.

Let θ_o be the initial parameters.

$\theta \leftarrow \theta_o$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

while $j < p$ **do**

 Update θ by running the training algorithm for n steps.

$i \leftarrow i + n$

$v' \leftarrow \text{ValidationSetError}(\theta)$

if $v' < v$ **then**

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$v \leftarrow v'$

else

$j \leftarrow j + 1$

end if

end while

Best parameters are θ^* , best number of training steps is i^*

Algorithm 7.2 A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

Let $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ be the training set.

Split $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ into $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$ and $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$ respectively.

Run early stopping (Algorithm 7.1) starting from random $\boldsymbol{\theta}$ using $\mathbf{X}^{(\text{subtrain})}$ and $\mathbf{y}^{(\text{subtrain})}$ for training data and $\mathbf{X}^{(\text{valid})}$ and $\mathbf{y}^{(\text{valid})}$ for validation data. This returns i^* , the optimal number of steps.

Set $\boldsymbol{\theta}$ to random values again.

Train on $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ for i^* steps.

Zaustavljanje učenja konačnog modela prema validacijskoj grešci

Algorithm 7.3 Meta-algorithm using early stopping to determine at what objective value we start to overfit, then continue training until that value is reached.

Let $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ be the training set.

Split $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ into $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$ and $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$ respectively.

Run early stopping (Algorithm 7.1) starting from random θ using $\mathbf{X}^{(\text{subtrain})}$ and $\mathbf{y}^{(\text{subtrain})}$ for training data and $\mathbf{X}^{(\text{valid})}$ and $\mathbf{y}^{(\text{valid})}$ for validation data. This updates θ .

$\epsilon \leftarrow J(\theta, \mathbf{X}^{(\text{subtrain})}, \mathbf{y}^{(\text{subtrain})})$

while $J(\theta, \mathbf{X}^{(\text{valid})}, \mathbf{y}^{(\text{valid})}) > \epsilon$ **do**

 Train on $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ for n steps.

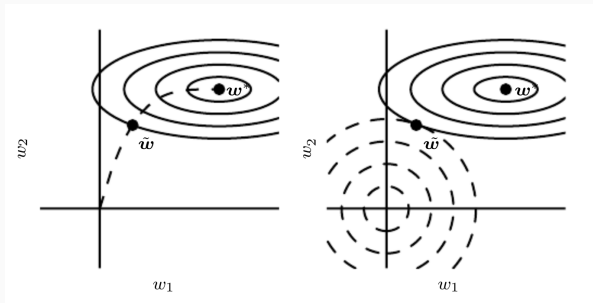
end while

Teoretska analiza ranog zaustavljanja

Pretpostavimo da prođemo kroz τ iteracija postupka gradijentnog spusta, s korakom ϵ , i da je gradijent u svakom koraku ograđen.

Tada je prostor parametara koji možemo dosegnuti iz početne vrijednosti parametara Θ_0 omeđen.

Veličina kugle ovisi o $\tau\epsilon$: što je ta vrijednost veća, to je kugla veća.



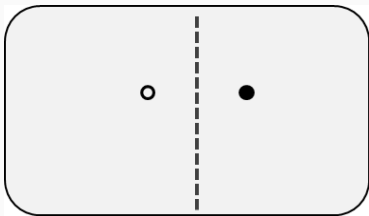
- Regularizacija
- Regularizacija norme vektora parametara modela
- Regularizacija generiranjem podataka i unošenjem šuma
- Regularizacija ranim zaustavljanjem
- Polunadzirano i višezadaćno učenje, dijeljenje parametara
- Regularizacija baggingom i dropout

polunadzirano učenje kao regularizacija

Polunadzirano učenje (eng. semi-supervised learning) je učenje $P(\mathbf{y}|\mathbf{x})$ koristeći i označene podatke (uzorke iz $P(\mathbf{x}, \mathbf{y})$) i neoznačene podatke (uzorke iz $P(\mathbf{x})$).

Možemo konstruirati model koji minimizira i nadzirani gubitak ($-\log P(\mathbf{y}|\mathbf{x})$) i nenadzirani gubitak (npr. $-\log P(\mathbf{x})$).

Dodatni signal za učenje može dovesti do boljeg semantičkog sadržaja dijeljenih značajki i poboljšati generalizaciju.

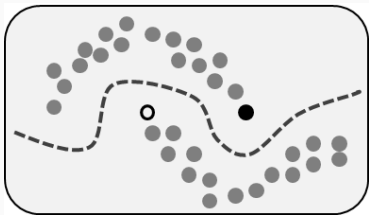


polunadzirano učenje kao regularizacija

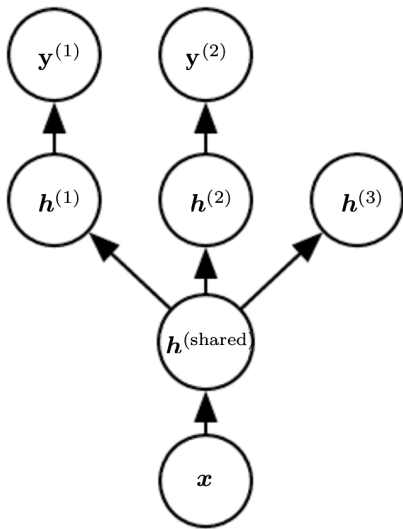
Polunadzirano učenje (eng. semi-supervised learning) je učenje $P(\mathbf{y}|\mathbf{x})$ koristeći i označene podatke (uzorke iz $P(\mathbf{x}, \mathbf{y})$) i neoznačene podatke (uzorke iz $P(\mathbf{x})$).

Možemo konstruirati model koji minimizira i nadzirani gubitak ($-\log P(\mathbf{y}|\mathbf{x})$) i nenadzirani gubitak (npr. $-\log P(\mathbf{x})$).

Dodatni signal za učenje može dovesti do boljeg semantičkog sadržaja dijeljenih značajki i poboljšati generalizaciju.



Višezadaćno učenje kao regularizacija



Npr:

$y^{(1)} \dots p(c|x)$

$y^{(2)} \dots p(x)$

$h^{(3)} \dots$ latentna reprezentacija
(za zamjenski gubitak)

Regularizacija vezanjem i dijeljenjem parametara

Vezanje parametara: ne znamo kakve parametre trebamo, ali znamo da postoji veza između dva problema

$$\Omega(\mathbf{w}^{(A)}, \mathbf{w}^{(B)}) = \|\mathbf{w}^{(A)} - \mathbf{w}^{(B)}\|_2^2$$

Dijeljenje parametara: u slučajevima kad želimo da skupovi parametara budu isti. Prednost u odnosu na vezanje parametara: trebamo čuvati samo jedan skup parametara.

Rijetke reprezentacije kao regularizator

- L_1 regularizacija težina \rightarrow rijetki modeli:

$$\begin{array}{ccc} \begin{bmatrix} 18 \\ 5 \\ 15 \\ -9 \\ -3 \end{bmatrix} & = & \begin{bmatrix} 4 & 0 & 0 & -2 & 0 & 0 \\ 0 & 0 & -1 & 0 & 3 & 0 \\ 0 & 5 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & -4 \\ 1 & 0 & 0 & 0 & -5 & 0 \end{bmatrix} \begin{bmatrix} 2 \\ 3 \\ -2 \\ -5 \\ 1 \\ 4 \end{bmatrix} \\ \mathbf{y} \in \mathbb{R}^m & & \mathbf{A} \in \mathbb{R}^{m \times n} \quad \mathbf{x} \in \mathbb{R}^n \end{array}$$

$$\tilde{J}(\Theta; \mathbf{X}, \mathbf{y}) = J(\Theta; \mathbf{X}, \mathbf{y}) + \alpha \|\Theta\|_1$$

Rijetke reprezentacije kao regularizator

- L_1 regularizacija aktivacija \rightarrow rijetke reprezentacije:

$$\begin{array}{ccc} \begin{bmatrix} -14 \\ 1 \\ 19 \\ 2 \\ 23 \end{bmatrix} & = & \begin{bmatrix} 3 & -1 & 2 & -5 & 4 & 1 \\ 4 & 2 & -3 & -1 & 1 & 3 \\ -1 & 5 & 4 & 2 & -3 & -2 \\ 3 & 1 & 2 & -3 & 0 & -3 \\ -5 & 4 & -2 & 2 & -5 & -1 \end{bmatrix} \begin{bmatrix} 0 \\ 2 \\ 0 \\ 0 \\ -3 \\ 0 \end{bmatrix} \\ \mathbf{y} \in \mathbb{R}^m & & \mathbf{B} \in \mathbb{R}^{m \times n} \quad \mathbf{h} \in \mathbb{R}^n \end{array}$$

$$\tilde{J}(\Theta; \mathbf{X}, \mathbf{y}) = J(\Theta; \mathbf{X}, \mathbf{y}) + \alpha \|\mathbf{h}\|_1$$

- Regularizacija
- Regularizacija norme vektora parametara modela
- Regularizacija generiranjem podataka i unošenjem šuma
- Regularizacija ranim zaustavljanjem
- Polunadzirano i višezadaćno učenje, dijeljenje parametara
- Regularizacija baggingom i dropout

Bagging

Bagging (ili **b**ootstrap **a**ggregating) [breiman96ml]: jednostavna ali efikasna metoda za poboljšanje generalizacijske točnosti

Neka je zadan skup za učenje $\mathbf{X} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$.

Iz \mathbf{X} uzorkujemo M podskupova $\mathbf{X}_m = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{N'}$

- uzorci se mogu ponavljati (uzorkovanje sa zamjenom)

Na svakom uzorku učimo jednu instancu modela: $f(\mathbf{x}, \Theta_m)$, te iz dobivenih instanci formiramo **ansambl**

Zaključivanje provodimo usrednjavanjem predikcija članova ansambla:

$$\hat{\mathbf{y}} = \frac{1}{M} \sum_{m=1}^M f(\mathbf{x}, \Theta_m)$$

Regularizacijski efekt bagginga

Pretpostavimo da smo naučili M skalarnih regresijskih modela:

- opišimo pogrešku i -tog modela varijablom ϵ_i
- pogreška ansambla regresora tada je: $E_a = \frac{1}{M} \sum_i \epsilon_i$

Pretpostavimo zajedničku distribuciju grešaka modela: $\epsilon \sim \mathcal{N}(\mathbf{0}, \Sigma)$

- dijagonalni elementi kovarijacijske matrice: $\mathbb{E}[\epsilon_i^2] = \sigma_{i,i} = v$,
- ne-dijagonalni elementi su: $\mathbb{E}[\epsilon_i \epsilon_j] = \sigma_{i,j} = c$.

Regularizacijski efekt bagginga (2)

Očekivanje kvadrata pogreške ansambla E_a^2 tada je:

$$\begin{aligned}\mathbb{E} \left[\left(\frac{1}{M} \sum_i \epsilon_i \right)^2 \right] &= \frac{1}{M^2} \mathbb{E} \left[\sum_i \left(\epsilon_i^2 + \sum_{i \neq j} \epsilon_i \epsilon_j \right) \right] \\ &= \frac{1}{M} v + \frac{M-1}{M} c\end{aligned}$$

Dva ekstremna slučaja:

1. Greške pojedinih regresora su potpuno korelirane $v = c$: ansambl tada griješi kao i pojedini modeli.
2. Greške pojedinih regresora su potpuno nekorelirane $c = 0$: greška ansambla tada opada linearno s brojem regresora M .

Bagging: ilustracija

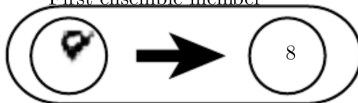
Original dataset



First resampled dataset



First ensemble member



Second resampled dataset



Second ensemble member



Usrednjavanje modela

Bitno je da modeli daju drugačije predikcije, učenje iz raznih podskupova je samo jedan način da se to omogući.

U dubokim modelima različitost možemo postići:

- različitim hiper-parametrima modela i postupka učenja
- različitim slučajnim inicijalizacijama
- različitim rasporedom podataka u mini-grupama

Ti efekti su dovoljni da naučeni modeli rade različite greške.

Usrednjavanje modela u praksi radi jako dobro: to je način kako se dobivaju najbolji rezultati (i pobjeđuje na natjecanjima).

Tehnika regularizacije dubokih modela:

- srodna ansambliranju, računski manje zahtjevna

Intuicija: efikasni ansambl velikog broja različitih modela.

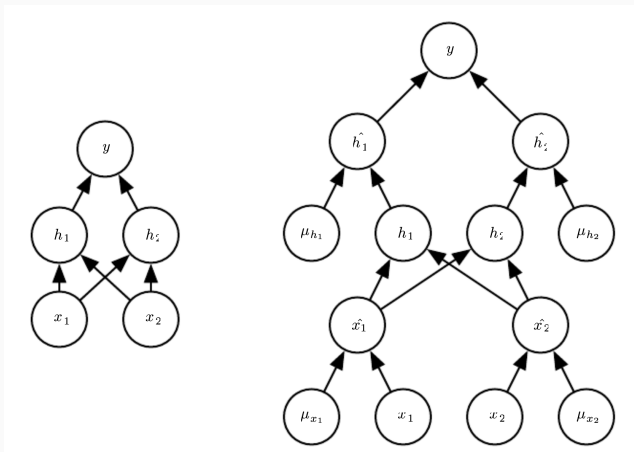
Kako dobiti veliki broj različitih modela?

- iz početnog modela formirati podmodele gašenjem/uklanjanjem aktivacija: latentnih i ulaznih značajki
- u praksi, odabir provodimo množenjem vektora aktivacija sloja \mathbf{x} s binarnom slučajnom maskom $\boldsymbol{\mu} \sim B(1, p)^N$:

$$\mathbf{x}' = \mathbf{x} \odot \boldsymbol{\mu}$$

- tipično, za ulazne značajke $p = 0.8$, a za skrivene značajke $p = 0.5$

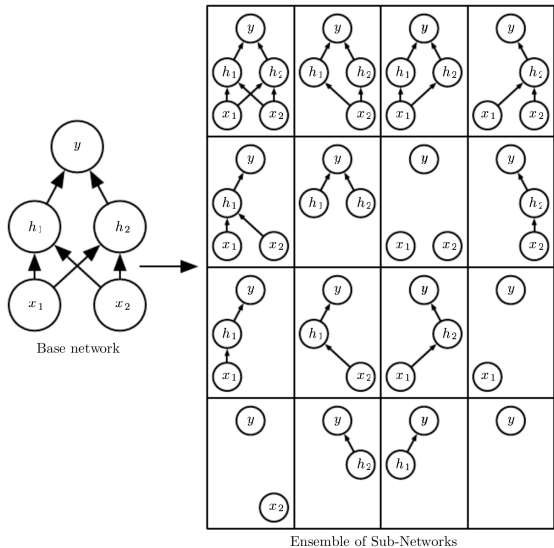
Dropout kao ansambl 2^N modela



$$\boldsymbol{\mu} = [\mu_{x_1}, \mu_{x_2}, \mu_{h_1}, \mu_{h_2}]$$

$$p(\boldsymbol{\mu}) = p(\mu_{x_1})p(\mu_{x_2})p(\mu_{h_1})p(\mu_{h_2})$$

Dropout kao ansambl 2^N modela (2)

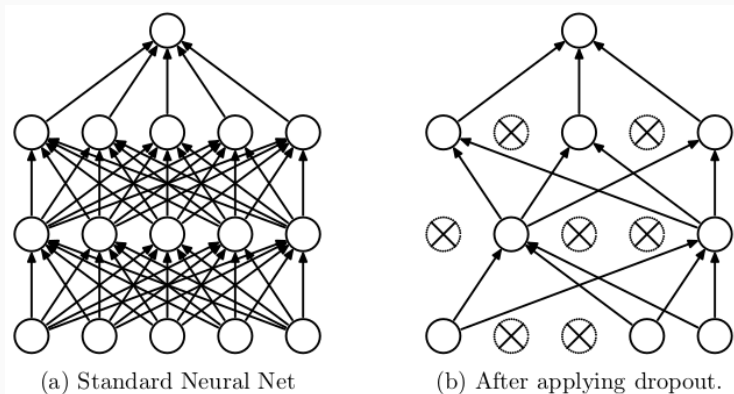


U praksi, $N > 100$: većina podmodela ima put od ulaza do izlaza

Dropout u fazi učenja

Dropout pristup: zasebno učiti dijelove istog modela

- formalno, minimiziramo: $\mathbb{E}_{\mu} J(\Theta, \mu) = \sum_{\mu} J(\Theta, \mu) p(\mu)$
- eksponencijalnu složenost možemo primiriti uzorkovanjem μ u svakoj iteraciji učenja (nepristrana procjena gradijenta)



Dropout u fazi ispitivanja

U slučaju bagginga, predikcija ansambla za podatak \mathbf{x} je:

$$\frac{1}{M} \sum_m p(\mathbf{y}|\mathbf{x}, \Theta_m)$$

Kod dropouta ansambliramo podmodele:

$$p(\mathbf{y}|\mathbf{x}, \Theta) = \sum_{\mu} p(\mu) p(\mathbf{y}|\mathbf{x}, \mu) \quad (1)$$

- prikazani izraz računa **očekivanje** preko distribucije maski $p(\mu)$ koje definiraju podmodele
- broj članova u izrazu je 2^N gdje je N broj aktivacija.

Dropout u fazi ispitivanja: uzorkovanje pod-modela

U praksi možemo raditi s M podmodela ($M \ll 2^N$)

Uzorkujemo $\boldsymbol{\mu}_m \sim p(\boldsymbol{\mu})$ i određujemo izlaze podmodela $m \in 1..M$:

$$\hat{\mathbf{y}}_m = p(\mathbf{y}|\mathbf{x}, \boldsymbol{\mu}_m)$$

Konačna predikcija je srednja vrijednost izlaza podmodela:

$$\hat{\mathbf{y}} = \frac{1}{M} \sum_{m=1}^M \hat{\mathbf{y}}_m$$

- problem: za ovo trebamo M unaprijednih prolaza...
- prednost: možemo izračunati varijancu predikcije...

Dropout u fazi ispitivanja: skaliranje težina

Ako aritmetičku sredinu zamjenimo geometrijskom:

$$\tilde{p}(\mathbf{y}|\mathbf{x}, \Theta) = \left(\prod_{\mu} p(\mathbf{y}|\mathbf{x}, \Theta_{\mu}) \right)^{2^{-n}} \dots$$

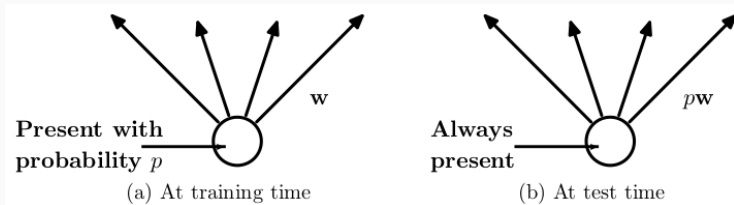
... te re-normaliziramo predikcije na distribuciju preko izlaza:

$$\hat{p}(\mathbf{y}|\mathbf{x}, \Theta) = \frac{\tilde{p}(\mathbf{y}|\mathbf{x}, \Theta)}{\sum_{c=1}^C \tilde{p}(y_c = 1|\mathbf{x}, \Theta)}$$

Tada dropout u fazi ispitivanja možemo aproksimirati **skaliranjem težina** modela.

Dropout u fazi ispitivanja

Izlazne težine aktivacije (= stupci matrice \mathbf{W}) množimo s vjerojatnošću njene prisutnosti tijekom učenja: $\mathbf{w}'_{h_i} = p(\mu_{h_i})\mathbf{w}_{h_i}$



[srivastava15jmlr]

- intuicija: osigurati da očekivani ulaz čvora u fazi ispitivanja bude jednak očekivanom ulazu čvora u fazi učenja
- prednost: jedan umesto M unaprijednih prolaza.

Druga mogućnost: skalirati aktivacije **prilikom učenja**: $h'_i = h_i/p(\mu_{h_i})$

Dropout: prednosti

Moguća primjena na razne tipove modela (konvolucijski, povratni) bez modifikacije funkcije cilja.

Jednostavna kombinacija s ostalim metodama regularizacije.

Računski znatno manje zahtjevan od bagginga

U usporedbi s ne-regulariziranim modelom:

- složenije učenje (više iteracija, npr. $2\times$)
- jednako brzo zaključivanje

Dropout: kada?

Dropout efektivno smanjuje kapacitet modela; najbolji rezultati postižu se povećanjem modela i duljim treniranjem. Za jako velike skupove podataka nije praktičan.

Kod konvolucijskih modela mogu se izbacivati:

- pojedinačne aktivacije
- cijele mape značajki (`drop_channel`)

Normalizacija nad grupom (eng. batch norm) također ima i regularizacijski efekt: dodavanje aditivnog i multiplikativnog šuma koji ovisi o podacima u mini-grupi. Popularna alternativa dropoutu.

Primjer: skaliranje težina u višerazrednoj logističkoj regresiji

Nema skrivenih slojeva \Rightarrow dropout samo na ulaznim značajkama \mathbf{x} .

Za višerazrednu logističku regresiju: geometrijska sredina predikcija eksponencijalnog broja pod-modela **ekvivalentna skaliranju težina**

Višerazredna logistička regresija:

$$p(\mathbf{y} = y | \mathbf{x}, \Theta) = \text{softmax}(\mathbf{W} \cdot \mathbf{x} + \mathbf{b})_y$$

Razred pod-modela definiran binarnim vektorom μ :

$$p(\mathbf{y} = y | \mathbf{x}, \Theta_\mu) = \text{softmax}(\mathbf{W} \cdot (\mu \odot \mathbf{x}) + \mathbf{b})_y$$

Primjer: skaliranje težina u višerazrednoj logističkoj regresiji

Dokažimo da predikciju eksponencijalno velikog dropout-ansambla uz $p=0.5$ uistinu možemo dobiti skaliranjem težina!

$$\begin{aligned}\tilde{p}(\mathbf{y} = y | \mathbf{x}, \Theta) &= \left(\prod_{\boldsymbol{\mu} \in \{0,1\}^n} p(\mathbf{y} = y | \mathbf{x}, \Theta_{\boldsymbol{\mu}}) \right)^{2^{-n}} \\&= \frac{\left(\prod_{\boldsymbol{\mu} \in \{0,1\}^n} \exp(\mathbf{W}_{y,:}(\boldsymbol{\mu} \odot \mathbf{x}) + b_y) \right)^{2^{-n}}}{\left(\prod_{\boldsymbol{\mu} \in \{0,1\}^n} \sum_{y'} \exp(\mathbf{W}_{y',:}(\boldsymbol{\mu} \odot \mathbf{x}) + b_{y'}) \right)^{2^{-n}}} \\&= C \cdot \left(\prod_{\boldsymbol{\mu} \in \{0,1\}^n} \exp(\mathbf{W}_{y,:}(\boldsymbol{\mu} \odot \mathbf{x}) + b_y) \right)^{2^{-n}} \\&= C \cdot \exp \left(2^{-n} \sum_{\boldsymbol{\mu} \in \{0,1\}^n} \mathbf{W}_{y,:}(\boldsymbol{\mu} \odot \mathbf{x}) + b_y \right)\end{aligned}$$

Primjer: skaliranje težina u višerazrednoj logističkoj regresiji

$$\begin{aligned}\tilde{p}(\mathbf{y} = y | \mathbf{x}, \Theta) &= C \cdot \exp \left(2^{-n} \sum_{\boldsymbol{\mu} \in \{0,1\}^n} \mathbf{W}_{y,:} (\boldsymbol{\mu} \odot \mathbf{x}) + b_y \right) \\ &= C \cdot \exp \left(\mathbf{W}_{y,:} \left(2^{-n} \sum_{\boldsymbol{\mu} \in \{0,1\}^n} \boldsymbol{\mu} \right) \odot \mathbf{x} + 2^{-n} \sum_{\boldsymbol{\mu} \in \{0,1\}^n} b_y \right) \\ &= C \cdot \exp \left(\mathbf{W}_{y,:} (\mathbf{0.5} \odot \mathbf{x}) + b_y \right) \\ &= C \cdot \exp \left((\mathbf{W}_{y,:} \odot \mathbf{0.5}) \cdot \mathbf{x} + b_y \right)\end{aligned}$$

Lako je pokazati da bi se drugi izbori $p(\boldsymbol{\mu})$ postigli drugačijim iteriranjem kroz podmodele. Tada bismo imali:

$$\tilde{p}(\mathbf{y} = y | \mathbf{x}, \Theta) = C \cdot \exp \left(\left(\mathbf{W}^\top \odot p(\boldsymbol{\mu}) \right)_{y,:}^\top \mathbf{x} + b_y \right)$$

Zadatak

Pretpostavke:

- 3-dimenzionalni ulaz: $\mathbf{x} \in \mathbb{R}^3$.
- dvije klase na izlazu $\mathbf{y} \in \mathbb{R}^2$.
- model: višerazredna logistička regresija

$$\mathbf{y} = \text{softmax}(\mathbf{W}\mathbf{x} + \mathbf{b})$$

- ulaz modela smo regularizirali dropoutom (bez mijenjanja aktivacija) uz $p(\mu_{x_1}) = 0.2$, $p(\mu_{x_2}) = 0.5$, $p(\mu_{x_3}) = 0.7$.

Zadatak: izračunajmo ulaz u softmax (eng. logits) za podatak $\mathbf{x} = [1, 1, 1]$, ako su parametri modela nakon učenja bili:

$$\mathbf{W} = \begin{bmatrix} -0.1 & 0.4 & -0.6 \\ 0.2 & -0.3 & 0.5 \end{bmatrix} \quad \mathbf{b} = \begin{bmatrix} 0.2 \\ -0.2 \end{bmatrix}$$