

SPR day 9

Probabilistic Generalized Linear Models

Reading: Bishop PRML Sec 4.3

- Outline:
- (1) Discriminative vs. Generative
 - (2) Generalized Linear Models
 - (3) Sigmoid Function
 - (4) 2-class Probabilistic Logistic Regression
 - ML estimation and MAP estimation
 - Review possible strategies to solve $\min_{w \in \mathbb{R}^M} d(w)$
 - (5) Understanding 2nd order gradient methods

Discriminative vs Generative Models

Consider supervised learning task:

Given training data $\{x_n, y_n\}_{n=1}^N$

features $x_n \in \mathbb{R}^D$
labels $y_n \in \mathcal{Y}$

Learn to predict $p(y_* | x_*)$ labels-given-features

Two Approaches to probabilistic modeling

Discriminative (Sec 4.3)

Generative (Sec 4.2)

Model:
$$p(y|x) = \prod_{n=1}^N P_w(y_n | x_n)$$

only y
is r.v.

x treated as fixed/known

Parameters: w generates label given feature

Prediction:
$$p(y_* | x_*, w)$$

directly use
likelihood

Training:
$$\max_w \prod_n p(y_n | x_n, w)$$

Pro: - simpler. fewer parameters.
- directly solve supervised task

Con: - cannot predict if x has missing values

Model:
$$p(x, y) = \prod_{n=1}^N P_\theta(x_n | y_n) P_\pi(y_n)$$

both
 x & y
as r.v.

Parameters: θ generate x given y
 π generates y

Prediction:
$$p(y_* | x_*) = \frac{P_\theta(x_* | y_*) P_\pi(y_*)}{\sum_{y' \in \mathcal{Y}} P_\theta(x_* | y') P_\pi(y')}$$

via
Bayes Rule

Training:
$$\arg \max_{\theta, \pi} \sum_n \log P_\theta(x_n | y_n) + \log P_\pi(y_n)$$

Pro: - can predict if x w/ missing values
- can sample new x from $p(x|y)$

Con: - more complex prediction
- more parameters

Generalized Linear Models

Extend linear regression to tasks with constrained output space
(output space: real line)

Training Data: $\{x_n, t_n\}_{n=1}^N$ where $\phi(x_n) \in \mathbb{R}^M$, $t_n \in \mathcal{Y}$
output space

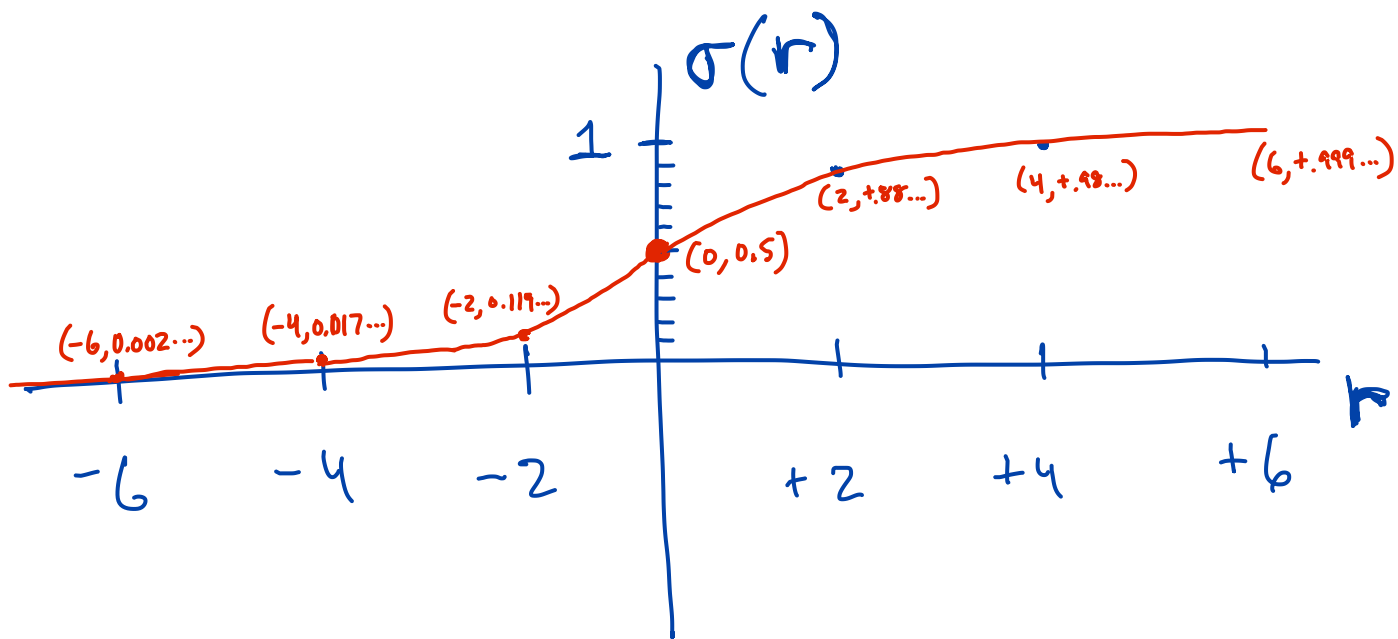
Model: $p(t_{1:N} | x_{1:N}) = \prod_{n=1}^N Q(t_n | f(w^T \phi(x_n)))$
chosen PMF or PDF over output space \mathcal{Y}
link function or activation function Determines mean $\rightarrow E[t_n] = f(w^T \phi(x_n))$
Key Idea

Output Space \mathcal{Y}	name	distribution Q	link function f
real $(-\infty, +\infty)$	linear regression	Normal	$f(w, x) = w^T \phi(x)$
integers $\{0, 1, 2, \dots\}$	Poisson regression	Poisson	$f(w, x) = e^{w^T \phi(x)}$
binary $\{0, 1\}$	Logistic regression	Bernoulli	$f(w, x) = \sigma(w^T \phi(x))$ <small>logistic sigmoid</small>
	Probit regression	Bernoulli	$f(w, x) = \Phi(w^T \phi(x))$ <small>Normal CDF</small>
C-ary classification <small>0 0 0 ... 1 0 0 0 ... 10 ! 1 0 0 ... 00</small>	Multi-class logistic regression	Categorical	$f(w, x) = \text{Softmax}(W^T \phi(x))$ <small>$W: M \times C$ matrix <small>vector in \mathbb{R}^C</small></small>
C levels w/ known ranking <small>e.g. ratings 0-5 stars e.g. cold, warm, hot</small>	Ordinal regression	Ordinal	<small>look up if interested!</small>
positive reals $(0, +\infty)$	Exponential regression	Exponential	

Sigmoid function

Let $\sigma(r)$ denote logistic sigmoid function

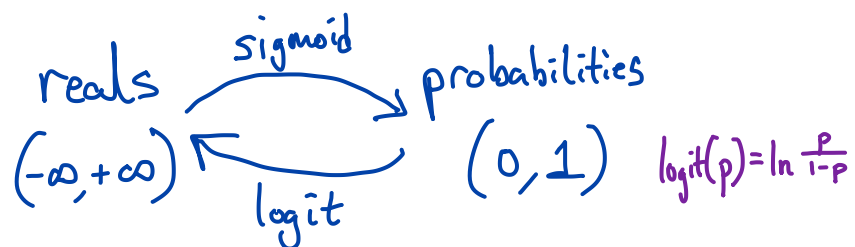
$$\sigma(r) = \frac{1}{1+e^{-r}} = \frac{e^r}{e^r+1}$$



Properties:

- Maps real line to probability interval (0, 1)
- Monotonic: $r_B > r_A$ implies $\sigma(r_B) > \sigma(r_A)$

- Invertible:



Probabilistic Logistic Regression

Prior: $p(w) = N(w \mid m_0, S_0)$
 $= N(w \mid \vec{0}, \alpha^{-1} I_M)$ usually

Likelihood $p(t_n \mid w) = \text{BernPMF}(t_n \mid \sigma(w^T \phi(x_n)))$

Link function is logistic sigmoid
Monotonic transform from
 $(-\infty, +\infty)$ to $(0, 1)$

Point Estimation Strategies

ML estimate: $\underset{w \in \mathbb{R}^M}{\text{argmax}} \sum_{n=1}^N \log p(t_n \mid w)$

MAP estimate: $\underset{w \in \mathbb{R}^M}{\text{argmax}} \sum_{n=1}^N \log p(t_n \mid w) + \log p(w)$

Next Time: Posterior Estimation: $p(w \mid t_{1:N})$

Strategies for Optimization

$$\underset{w \in \mathbb{R}^M}{\operatorname{argmin}} d(w)$$

Convention: always minimize

ML/MAP for
Linear Regression

ML/MAP for
Logistic Reg.

Analytical methods

- Set up $\nabla_w d = 0$
- Manipulate to find $w^* = \dots$
closed form expression

rarely works
but fast + exact
once implemented



No closed form solution exists.

Gradient methods

stepsize $\epsilon_t > 0$
gradient $g \in \mathbb{R}^M$
Hessian $H \in \mathbb{R}^{M \times M}$
sym
p.d.

- 1st order gradient descent
 $w_{t+1} \leftarrow w_t - \epsilon_t g(w_t)$

Many iterations,
each cheap:
 $O(M)$



- 2nd order gradient descent
 $w_{t+1} \leftarrow w_t - \epsilon_t H(w_t)^{-1} g(w_t)$

Few iters,
each expensive:
 $O(M^3)$
inverting matrix



Other methods

- grid search
- random search
- Nelder-Mead
- ... many others possible

usually very
inefficient in
high dimensions
($M > 3$)



ML/MAP estimation: Gradients + Hessians

For Logistic Regression: $d(w) = - \sum_{n=1}^N \log \text{BernPMF}(t_n | \sigma(w^T \phi(x_n)))$
 - $\log \text{NormPDF}(w | 0, \alpha^{-1} I_M)$

gradient $\nabla_w d$
 shape: $(M, 1)$

Hessian $\nabla_w \nabla_w d$
 shape: (M, M)

Linear
 Regr.

$$\beta \Phi^T (\underbrace{\Phi w}_{\text{predicted}} - \underbrace{t}_{\text{true label}}) + \alpha w$$

$$\beta \Phi^T \Phi + \alpha I_M$$

Logistic
 Regr.

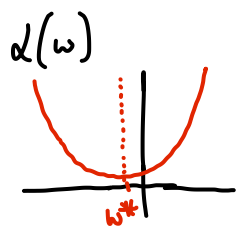
$$\Phi^T (\underbrace{\sigma(\Phi w)}_{\text{prediction (probability)}} - \underbrace{t}_{\text{true binary label}}) + \alpha w$$

$$\Phi^T R(w) \Phi + \alpha I_M$$

where $R(w) = \begin{bmatrix} r_1(w) & & \\ & \ddots & \\ & & r_N(w) \end{bmatrix}$ diagonal $N \times N$ matrix
 All diagonal entries $r_i > 0$
 Off diagonal entries = 0

$$r_n(w) = \sigma(w^T \phi(x_n)) \sigma(-w^T \phi(x_n))$$

always > 0 since $\sigma(a) > 0$ for all a



Unique minima w^* exists when objective is convex.
 Second derivative test is positive.
 In M -dims, Hessian is positive definite.

Question:

Is Linear Regression MAP convex?

Yes. Hessian = $\underbrace{\alpha I_M}_{\text{pos. def. if } \alpha > 0} + \underbrace{\Phi^T \Phi}_{\text{p.s.d.}}$. Sum of p.d. and psd is positive definite.

Is Logistic Regression MAP convex?

Yes. R is a positive diagonal, so $\Phi^T R \Phi$ still psd, $\alpha I + \Phi^T R \Phi$ is p.d.

2nd order gradient methods

Solving an optimization problem:

$$\arg \min_{w \in \mathbb{R}^M} \alpha(w)$$

Let $g(w) = \nabla_w \alpha$

shape

$$(M, 1)$$

$$H(w) = \nabla_w \nabla_w \alpha$$

$$(M, M)$$

1st order GD

$$w^{\text{new}} \leftarrow w^{\text{old}} - \epsilon g(w^{\text{old}})$$

if $M=1$ special case:

$$w^{\text{new}} \leftarrow w^{\text{old}} - \epsilon \cdot \alpha'(w^{\text{old}})$$

2nd order GD

$$w^{\text{new}} \leftarrow w^{\text{old}} - \epsilon H(w^{\text{old}})^{-1} g(w^{\text{old}})$$

$$w^{\text{new}} \leftarrow w^{\text{old}} - \epsilon \cdot \frac{\alpha'(w^{\text{old}})}{\alpha''(w^{\text{old}})}$$

looks like Newton's method!

Linear Regression: 2nd order GD

Using formulas from previous page:

$$w^{\text{new}} = w^{\text{old}} - \epsilon [\beta \Phi^T \Phi + \alpha I_M]^{-1} (\beta \Phi^T \Phi w^{\text{old}} - \beta \Phi^T z + \alpha w^{\text{old}})$$

$$= w^{\text{old}} - \epsilon (\beta \Phi^T \Phi + \alpha I_M)^{-1} [(\beta \Phi^T \Phi + \alpha I_M) w^{\text{old}} - \beta \Phi^T z]$$

$$= w^{\text{old}} - \epsilon w^{\text{old}} + \epsilon (\beta \Phi^T \Phi + \alpha I_M)^{-1} \beta \Phi^T z$$

If stepsize $\epsilon=1$, we get:

$$w^{\text{new}} = (\Phi^T \Phi + \frac{\kappa}{\beta} I_M)^{-1} \Phi^T z$$

optimal MAP estimate in one step!

Anytime loss is quadratic, 2nd order gradient update with step size $\epsilon=1$ will find global minima in one step.

Logistic Regression w/ 2nd order GD

Still gold standard for ML/MAP estimation,
but requires many iterations (unlike linear regression ML estimation using 2nd order methods).

Update looks like this (assuming ML estimation. MAP is similar.)

$$w^{\text{new}} \leftarrow w^{\text{old}} - \epsilon \left(\Phi^T R(w^{\text{old}}) \Phi \right)^{-1} \Phi^T \left(\underbrace{\nabla(\Phi w^{\text{old}})}_{\hat{y}} - t \right)$$

Assume $\epsilon=1$ and rearranging terms

$$\leftarrow \left(\Phi^T R \Phi \right)^{-1} \left(\Phi^T R \Phi w^{\text{old}} - \Phi^T \hat{y} + \Phi^T t \right)$$

$$\leftarrow \left(\Phi^T R \Phi \right)^{-1} \left(\Phi^T R z \right)$$

where $z \in \mathbb{R}^N$
 $z = \Phi w^{\text{old}} - R^{-1}(\hat{y} - t)$

Interpret as "least squares"

with per-example "weights" r , $R = \text{diag}(r)$

and outcomes z , both of which depend on previous value of w .

Thus, the name

"Iteratively Re-weighted Least Squares" (IRLS)

Often used for Generalized Linear Models
to do ML estimation of weights $w \in \mathbb{R}^M$

while not converged:

update R given w

update z given w

update $w \leftarrow \left(\Phi^T R \Phi \right)^{-1} \Phi^T R z$