1. (a).
   i. B     ii. B,C

(b)
   i. CD     ii. A     iii. C
   iv. C     v. B     vi. B

(c)
   i. C,F     ii. BD     iii. AE

(d).
   i. BD     ii. BD     iii. CE

(e)
A: False. Logistic Regression cannot converge when data is linearly separable.

B: False. more leaves means more complicated hypothesis which will increase variance but decrease bias.

C:

D: True It dependends on support vectors which is a small subset of data

E:

F: True All the linear + identity activation layers can just be replaced with one linear function $a = W^T x + b$, applied on $\hat{y} = \text{sigmoid}(w^T x + b)$ which is simply a linear classifier.

G: False. Kernelized Perception doesn't maximize the margin, it only separates the data with a boundary.

H: True $p(x;\theta) = \frac{\pi}{\theta} \left(\frac{x}{\theta}\right)^{\pi-1} \exp\left\{-\left(\frac{x}{\theta}\right)^{\pi}\right\}$

$= \pi x^{\pi-1} \theta^{-\pi} \exp\left(-\theta^{-\pi} x^{\pi}\right)$

$= \pi x^{\pi-1} \exp\left(-\theta^{-\pi} x^{\pi} - \pi \log \theta\right)$,

$b(y): \pi x^{\pi-1}$

$\eta : -\theta^{-\pi}$

$a(\eta): \pi \log(\theta)$

$T(y): x^{\pi}$

(f):

A. YES        B. ~~NO~~ NO    C. NO    D. YES
E. NO        F. YES        G. YES        H. ~~NO~~ NO
I. YES        J. NO

## 2. CNN

(a): i, $m = n - d + 1$

ii, $W$ is a $(n-d+1) \times n$ matrix in the form:

$$W = \begin{bmatrix} \theta_1 \, \theta_2 \cdots \theta_d & & & \\ & \theta_1 \, \theta_2 \cdots \theta_d & & \\ & & \theta_1 \, \theta_2 \cdots \theta_d & \\ & & \cdots \cdots & \\ & & & \theta_1 \, \theta_2 \cdots \theta_d \end{bmatrix}$$

each row is a vector $(\theta_1, \theta_2 \cdots \theta_d)$ shifted right by (row−1) position.

iii, only $d$ parameters needed in the case without intercept, its fully connected equivalent on the other hand will need $(n-d+1) \cdot n$ parameters, ~~but only~~ ~~~~. Computationally since $d \ll n$, a large portion of $W$ will be $0$ which means forward/backward propagation will be much faster than a fully connect network. For both f/b propagation, CNN will only take $(n-d+1) \cdot d / (n-d+1) \cdot n = d/n$ of what a fully connect network needs in ~~comput~~ computation.

(b)

i: $\theta^{[1]}$ is a $d$ dimensional vector, it is used to make the diagonal matrix as indicated in (a). i ;

$W^{[2]}$ is $(n-d+1)$ dimensional vector ;

$b^{[2]}$ is just a scalar.

ii. The parameters of the network are: $\theta^{[1]}$, $W^{[2]}$ and $b^{[2]}$.

$\theta^{[1]} \in \mathbb{R}^{d}$ ; we CNN in the setting, $Z^{[1]}$, $a^{[1]}$ are $\mathbb{R}^{(n-d+1)}$ ;

$W^{[2]} \in \mathbb{R}^{n-d+1}$, $b^{(2)} \in \mathbb{R}$ ; hence, $Z^{[2]}$, $a^{[2]}$ are $\mathbb{R}$ ;

First write all the equations in the network. for ease of description, use $W^{[1]}$ as the CNN matrix composed by $\theta^{[1]}$ ;

$W^{[1]}$ is just $\mathbb{R}^{(n-d+1) \times n}$ matrix: $W^{[1]} = \begin{bmatrix} - \theta^{[1]} - \\ - \theta^{[1]} - \\ \ddots \\ - \theta^{[1]} - \end{bmatrix}$

$L(\hat{y}, y) = -\left((1-y) \log(1-\hat{y}) + y \log(\hat{y})\right)$

$\hat{y} = a^{[2]} = \delta(z^{[2]})$

$z^{[2]} = W^{[2]} a^{[1]} + b^{[2]}$

$a^{[1]} = Relu(z^{[1]})$

$z^{[1]} = W^{[1]} x$ ;

$(W^{[2]})$

$\dfrac{\partial L}{\partial W^{[2]}} = -\dfrac{\partial}{\partial W^{[2]}}\left((1-y) \log(1-\hat{y}) + y \log \hat{y}\right)$

$= \left((1-y)/(1-\hat{y}) - y/\hat{y}\right) \cdot \dfrac{\partial \delta}{\partial z^{[2]}} \dfrac{\partial z^{[2]}}{\partial W^{[2]}}$

$= \left((1-y)/(1-\hat{y}) - y/\hat{y}\right) \cdot \hat{y}(1-\hat{y}) \cdot a^{[1]}$

$= \left(\hat{y}(1-y) - y(1-\hat{y})\right) a^{[1]}$

$= (a^{[2]} - y) a^{[1]}$ ; $\quad W^{[2]}, a^{[1]} \in \mathbb{R}^{n-d+1}$ ; $a^{[2]} \in \mathbb{R}$ ;

$(b^{[2]})$ :

$\dfrac{\partial L}{\partial b^{[2]}} = -\dfrac{\partial}{\partial b^{[2]}}\left((1-y) \log(1-\hat{y}) + y \log \hat{y}\right)$

$= \left((1-y)/(1-\hat{y}) + y/\hat{y}\right) \cdot \dfrac{\partial \delta}{\partial z^{[2]}} \cdot \dfrac{\partial z^{[2]}}{\partial b^{[2]}}$

$= a^{[2]} - y$ ; $\quad b^{[2]}, a^{[2]}$ both are scalars;

$(a^{[1]})$ : $\Rightarrow$ NEXT PAGE

$(\theta^{[1]})$:

$\frac{\partial L}{\partial W^{[1]}} = \frac{\partial L}{\partial \hat{x}^{[2]}} \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial W^{[1]}}$

$= (a^{[2]} - y) \cdot W^{[2]} \cdot \beta \circ x$,    "$\circ$" here is outer product.

where $\beta = \frac{\partial a^{[1]}}{\partial z^{[1]}}$ is just element wise derivative of $a^{[1]}_i$ w.r.t $z^{[1]}_i$, with the ith element 1 when $z^{[1]}_i > 0$, otherwise $0$;

The dimension of $\frac{\partial L}{\partial W^{[1]}}$:

in the expression $(a^{[2]} - y) \cdot W^{[2]} \cdot \beta \circ x$, $(a^{[2]} - y)$ is a scalar; $W^{[2]}$ is $\mathbb{R}^{n-d+1}$; $\beta$ is $\mathbb{R}^{n-d+1}$; $x$ is $\mathbb{R}^n$; we have $\frac{\partial L}{\partial W^{[1]}}$ is $\mathbb{R}^{(n-d+1) \times n}$, the same as $W^{[1]}$;

For simplicity, let $A = \frac{\partial L}{\partial W^{[1]}}$; the gradient w.r.t $\theta^{[1]}$ is just an $\mathbb{R}^d$ vector, it's ith element is given by:

$\left( \frac{\partial L}{\partial \theta^{[1]}} \right)_i = \sum_{j=1}^{n-d+1} A_{(j,i)}$; $1 \le i \le d$;

Update rule is:

$W^{[2]} := W^{[2]} - \alpha \frac{\partial L}{\partial W^{[2]}}$,    $\frac{\partial L}{\partial W^{[2]}} = (a^{[2]} - y) a^{[1]}$

$b^{[2]} := b^{[2]} - \alpha \frac{\partial L}{\partial b^{[2]}}$,    $\frac{\partial L}{\partial b^{[2]}} = a^{[2]} - y$

$\theta^{[1]} := \theta^{[1]} - \alpha \frac{\partial L}{\partial \theta^{[1]}}$,    $\left( \frac{\partial L}{\partial \theta^{[1]}} \right)_i = \sum_{j=1}^{n-d+1} A_{(j,i)}$,   $A_{(j,i)}$ denotes the jth row, ith column.

(c): Now for $L(\hat{y}, y) = \frac{1}{2}(\hat{y}-y)^2$: we chain rule at derivation

$\frac{\partial L}{\partial w^{[2]}} = (\hat{y}-y) \cdot \hat{y}(1-\hat{y}) \cdot a^{[1]}$

$\quad = (a^{[2]}-y) \cdot a^{[2]}(1-a^{[2]}) \cdot a^{[1]}$ ;

dimension of $\frac{\partial L}{\partial w^{[2]}} \in \mathbb{R}^{n-d-1}$ ;

$\frac{\partial L}{\partial b^{[2]}} = (\hat{y}-y) \cdot \hat{y}(1-\hat{y})$

$\quad = (a^{[2]}-y) a^{[2]}(1-a^{[2]})$ ;

$\frac{\partial L}{\partial b^{[2]}} \in \mathbb{R}$ ;

$\frac{\partial L}{\partial w^{[1]}} = \frac{\partial L}{\partial a^{[2]}} \cdot \frac{\partial a^{[2]}}{\partial z^{[2]}} \cdot \frac{\partial z^{[2]}}{\partial a^{[1]}} \frac{\partial a^{[1]}}{\partial z^{[1]}} \frac{\partial z^{[1]}}{\partial w^{[1]}}$

$\quad \begin{cases} (a^{[2]}-y) \cdot a^{[2]}(1-a^{[2]}) \cdot w^{[2]} \circ x \ ; & \text{when } z^{[1]} \geqslant 0 \\ 0 \ ; & \text{when } z^{[1]} < 0 \ ; \end{cases}$

$\quad = (a^{[2]}-y) a^{[2]}(1-a^{[2]}) \cdot w^{[2]} \cdot \beta \circ x$ ;

the above $\beta = \frac{\partial a^{[1]}}{\partial z^{[1]}}$ is just a $\mathbb{R}^{n-d+1}$ vector, with the $i$th element 1 if $z_i^{[1]} \geqslant 0$, otherwise $\beta_i = 0$ ; the last product "$\circ$" is an outer product.

$\frac{\partial L}{\partial w^{[1]}} \in \mathbb{R}^{(n-d+1) \times n}$ ; for simplicity, let $\frac{\partial L}{\partial w^{[1]}} = A$ ;

the $i$th element of $\theta^{[1]}$'s gradient is:

$\left(\frac{\partial L}{\partial \theta^{[1]}}\right)_i = \sum_{j=1}^{n-d+1} A_{(j,i)}$. $1 \leq j \leq n-d+1$ ; $1 \leq i \leq d$ , $A_{(j,i)}$ is the $j$th row and $i$th column.

update rules:

$w^{[2]} = w^{[2]} - \alpha \frac{\partial L}{\partial w^{[2]}}$ ; $\frac{\partial L}{\partial w^{[2]}} = (a^{[2]}-y) \cdot a^{[2]}(1-a^{[2]}) \cdot a^{[1]}$ ;

$b^{[2]} := b^{[2]} - \alpha \frac{\partial L}{\partial b^{[2]}}$ ; $\frac{\partial L}{\partial b^{[2]}} = (a^{[2]}-y) a^{[2]}(1-a^{[2]})$

$\theta^{[1]} := \theta^{[1]} - \alpha \frac{\partial L}{\partial \theta^{[1]}}$ ; $\left(\frac{\partial L}{\partial \theta^{[1]}}\right)_i = \sum_{j=1}^{n-d+1} A_{(j,i)}$.

$\quad A = \frac{\partial L}{\partial \theta^{[1]}}$ as described above ; $A_{(j,i)}$ denotes the $j$th row, $i$th column of matrix $A$ ;

3. Linearity of Multinomial Naive Bayes.

Naive Bayes model contains $2k+1$ parameters, $k$ is the size of vocab.

$y \sim \text{Bernoulli}(\theta)$,

$x_i | y=0 \sim \text{Bernoulli}(\theta_{i|y=0})$

$x_i | y=1 \sim \text{Bernoulli}(\theta_{i|y=1})$

All $\theta$, $\theta_{i|y=1}$ and $\theta_{i|y=0}$ are given in lecture notes.

Naive Bayes decision boundary is given by:

$p(y|x) = 0.5$, from lecture notes, we know that:

$= p(x|y=1) \, p(y=1) \, / \, (p(x|y=0) \, p(y=0) + p(x|y=1) \, p(y=1))$

$\Rightarrow p(x|y=1) \, p(y=1) = p(x|y=0) \, p(y=0)$

$\Rightarrow \theta_{x|y=1} \, \theta = \theta_{x|y=0}(1-\theta) \Rightarrow \dfrac{\theta_{x|y=1}}{\theta_{x|y=0}} = \dfrac{1-\theta}{\theta}$

take log on both side, we have:

$\log \dfrac{\theta_{x|y=1}}{\theta_{x|y=0}} + \log \dfrac{\theta}{1-\theta} = 0$

$p(y=1|x)$

$= p(x|y=1) \, p(y=1) \, / \, p(x)$

$= (\prod_{i=1}^{n} (p(x_i|y=1)) \, p(y=1) \, / \, ((\prod_{i=1}^{n} p(x_i|y=0)) p(y=0) + (\prod_{i=1}^{n} p(x_i|y=1)) p(y=1))$

Set it to $0.5$, we have:

$(\prod_{i=1}^{n} p(x_i|y=1)) \, p(y=1) = (\prod_{i=1}^{n} p(x_i|y=0)) \, p(y=0)$

Which implies: $\dfrac{\prod_{i=1}^{n} p(x_i|y=1)}{\prod_{i=1}^{n} p(x_i|y=0)} = \dfrac{p(y=0)}{p(y=1)}$

take log on both sides and move the right to left:

$\sum_{i=1}^{n} \log \dfrac{p(x_i|y=1)}{p(x_i|y=0)} + \log \dfrac{p(y=1)}{p(y=0)} = 0$,

the duplicated words can be just replaced by $\phi(x_i) \cdot \log \dfrac{p(x_i|y=1)}{p(x_i|y=0)}$,

So, the above equation can be written as:

$\phi(x)^T \cdot \log \dfrac{\theta_{x|y=1}}{\theta_{x|y=0}} + \log \dfrac{\theta}{1-\theta} = 0$,

With: $m = \log \dfrac{\theta_{x|y=1}}{\theta_{x|y=0}}$ is a $k$-dimensional vector;

$c = \log \dfrac{\theta}{1-\theta}$ is a scalar.

4. Kernels.

(a). i. $|\Sigma|^n$

ii. $\lambda^4$

iii. Suppose $\psi(x) = \phi(x)/\sqrt{K_{ser}(x,x)}$ ;

We have $K_{norm}(x,z) = \psi(x)^T \psi(z) = \dfrac{K_{ser}(x,z)}{\sqrt{K_{ser}(x,x)K_{ser}}}$

$= K_{ser}(x,z)/\sqrt{K_{ser}(x,x)K_{ser}(z,z)}$

From class we know that if a matrix can be written as inner product ~~tmt~~ of two features $\psi(x), \psi(z)$, it is a kernel.

iv. normalized kernel can avoid a feature component to be scaled very small for non-contiguous substrings.
normalization makes the ~~string~~ long string kernel relatively invariant to the length of the document.

v. $2\lambda^x + \lambda^6$

(b). Refer to the problem set solution, we can know that, if $K$ is a kernel, it's polynomial ~~is~~ is still a kernel. $P(K(x,z))$

i. $K_1(x,z) = \exp(K(x,z))$

$= \sum_{i=0}^{\infty} \dfrac{1}{i!}(K(x,z))^i$ is a polynomial of a kernel, so it is still a kernel.

ii. $K_2(x,z) = \exp\left(-\dfrac{\|x-z\|^2}{\sigma^2}\right)$

$= \dfrac{1}{\exp\left(\frac{\|x\|^2+\|z\|^2}{\sigma^2}\right)} \exp\left(\dfrac{2x^Tz}{\sigma^2}\right)$

~~$G_1(\|x\|^2+\|z\|^2)$ $G_1 = \exp(\|x\|^2+\|z\|^2)$, matrix~~

$= \dfrac{1}{\exp\left(\frac{\|x\|^2}{\sigma^2}\right)} \cdot \dfrac{1}{\exp\left(\frac{\|z\|^2}{\sigma^2}\right)} \cdot \exp\left(\dfrac{2x^Tz}{\sigma^2}\right)$.

Let ~~$G_1(x,z)$~~

$G_1(x,z) = \psi(x)\psi(z)$ where $\psi(x) = \dfrac{1}{\exp(\|x\|^2/\sigma^2)}$ a scalar,

~~we~~ we know that $G_1(x,z)$ is a kernel;

Let $G_2(x,z) = \exp\left(\dfrac{2x^Tz}{\sigma^2}\right) = \exp(K_2(x,z))$ where $K_2(x,z) = 2x^Tz/\sigma^2$ is a kernel, from i we know $G_2$ is also a kernel.

$K_2(x,z) = G_1(x,z)G_2(x,z)$, as proved in PS 4.e, $K_2(x,z)$ is also a kernel.

iii. From Gaussian integral we have:

$$\sqrt{\pi/a} = \int_{-\infty}^{\infty} \exp(-a(t-b)^2)\,dt,$$

which implies:

$$1 = \int_{-\infty}^{\infty} \frac{1}{\sqrt{\pi/a}} \exp(-a(t-b)^2)\,dt\,;$$

together with:

$$k(x, z) = \exp\left(-\tfrac{1}{2}(x-z)^2\right), \quad \text{set } a = \tfrac{1}{2},\ b = (x+z)\,.$$

$$k(x,z) = \exp\left(-\tfrac{1}{2}(x-z)^2\right) \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{2}(t+x+z)^2\right)\,dt$$

$$= \int_{-\infty}^{\infty} \exp\left(-\tfrac{1}{2}\left(x^2+z^2-2xz\right)\right) \cdot \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{2}(t^2+x^2+z^2+2xt+2zt+2xz)\right)\,dt$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{2}\left(t^2+2x^2+2z^2+2xt+2zt\right)\right)\,dt$$

$$= \int_{-\infty}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\tfrac{1}{4}\left((t+2x)^2 + (t+2z)^2\right)\right)\,dt$$

$$= \int_{-\infty}^{\infty} (2\pi)^{-\frac{1}{4}} \exp\left(-\left(\tfrac{t+2x}{2}\right)^2\right) \cdot (2\pi)^{-\frac{1}{4}} \exp\left(-\left(\tfrac{t+2z}{2}\right)^2\right)\,dt$$

This gives us the final feature mapping function:

$$\phi(x, t) = (2\pi)^{-\frac{1}{4}} \exp\left(-\left(\tfrac{t+2x}{2}\right)^2\right)$$

(C): Kernelizing k-means.

The i-class centroid is now an infinite-dimension parameter which cannot be represented by computer, updating centroid directly is changed to update the indentities in Centroids calculation.

After Kernelizing, $u_j$ is still a linear combination of $\phi(x^{(i)})$:

$$u_j = \sum_{i=1}^{m} r_i\, \phi(x^{(i)}), \quad \text{where } r_i = \frac{1}{N_j} \text{ for } c^{(i)} = j, \text{ otherwise } r_i = 0.$$

The norm is now:

$$\|\phi(x^{(i)}) - u_j\|^2 = \left( \phi(x^{(i)})^T \phi(x^{(i)}) - 2\phi(x^{(i)})^T u_j + u_j^T u_j \right)$$

Plug in linear combination of $u_j$:

$$\|\phi(x^{(i)}) - u_j\|^2 = K(x^{(i)}, x^{(i)}) - 2\phi(x^{(i)})^T \sum_{j=1}^{m} r_j\, \phi(x^{(j)})$$
$$+ \sum_{i=1}^{m}\sum_{j=1}^{m} r_i r_j\, \phi(x^{(i)})^T \phi(x^{(j)})$$
$$= K(x^{(i)}, x^{(i)}) - \sum_{j=1}^{m} 2 r_j\, K(x^{(i)}, x^{(j)})$$
$$+ \sum_{i=1}^{m}\sum_{j=1}^{m} r_i r_j\, K(x^{(i)}, x^{(j)})$$

So the update rule should be changed to:

$$\tilde{c}^{(i)} = \arg\min_{j} \left( K(x^{(i)}, x^{(i)}) - 2\sum_{R=1}^{m} 2 r_R K(x^{(i)}, x^{(R)}) \right.$$
$$\left. + \sum_{R=1}^{m}\sum_{l=1}^{m} r_R r_l\, K(x^{(R)}, x^{(l)}) \right),$$

## 5. Trees and Random Forests.

### (a) ~~Classifying Trees~~

**i.** First prove that Gini loss is strictly concave:

For the 2 class case,

$G(R_m) = P_{m_1}(1-P_{m_1}) + P_{m_2}(1-P_{m_2})$ ; since $P_{m_1} + P_{m_2} = 1$, we have:

$G(R_m) = 2P_{m_1}(1-P_{m_1}) = 2P_{m_1} - 2P_{m_1}^2$ ;

take the 2nd order derivative of $G$ w.r.t $P_{m}$:

$G'' = -2 < 0$, which means Gini is strictly concave.

Then Let's prove that the weighted Gini loss is less or equal than the parent Gini; for simplicity, $P$ is used for $P_{m_1}$ for parent Gini:

The parent Gini loss: $G(R) = 2p(1-p)$.

The weight Gini loss of the children:

$$G(R_1, R_2) = \min_{j, t} \frac{|R_1|}{|R_1|+|R_2|} L(R_1) + \frac{|R_2|}{|R_1|+|R_2|} L(R_2)$$

$$= \min_{j,t} q \cdot 2 P_{11}(1-P_{11}) + (1-q) \cdot 2 \cdot P_{21}(1-P_{21})$$

$$\leq 2q \left( \frac{1}{q} P_{11}(1-P_{11}) + (1-q) P_{21}(1-P_{21}) \right)$$

$P_{11}$ means the first child's ~~proportion of~~ class 1 samples proportion,
$P_{21}$ means the second child's class 2 examples proportion;
$P$ means the class 1 samples proportion of the whole dataset;
~~let~~ We have $P_{11}|R_1| + P_{21}|R_2| = P(|R_1|+|R_2|)$, $R_1$, $R_2$ denotes the first/second trees. it is obvious that ~~$P_{11}, P_{21}$ lie~~
$P$ lies in between $P_{11}$ and $P_{21}$.

According to Jensen's inequality and strict concavity, we have
$$G(R_1, R_2) \leq 2q P_{11}(1-P_{11}) + 2(1-q) P_{21}(1-P_{21})$$
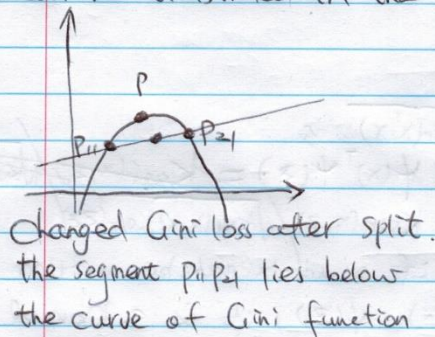$$\leq 2p(1-p)$$
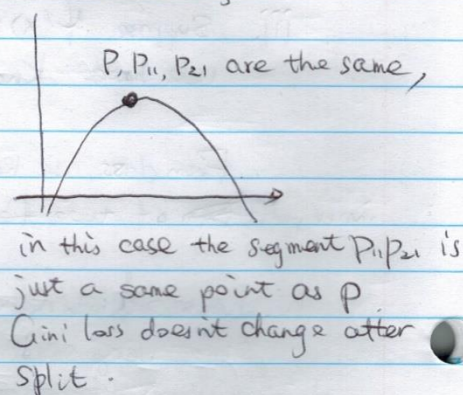$$= G(R),$$
Which is what need to be proved in i.


**ii.** in above explanation, we know that $P$ lies in between $P_{11}$ and $P_{21}$.
(in 2 class case, can be generalized to k classes)
if $P = P_{11} = P_{21}$, Which means the split doesn't change the class proportions in respective nodes, then ~~P~~ Gini Loss would not change after the split. in the strictly concave setting, because after split, class proportions of sub nodes doesn't change. ~~the~~ Next →

It can be illustrated in the diagram,; Left one is a Gini changed;
Right one is the case that
Gini doesn't change:



changed Gini loss after split.
the segment $P_{11} P_{21}$ lies below
the curve of Gini function

P, $P_{11}$, $P_{21}$ are the same,

in this case the segment $P_{11} P_{21}$ is
just a same point as P
Gini loss doesn't change after
split.

iii. In the misclassification case, apart from having all class proportions
to be the same, there is more cases where we could have the Gini
loss unchanged after split.
Suppose the max proportion of a class are all the same in ~~sub trees and~~
each of the children and in parent, the Gini loss will keep the same
after the split.
Suppose $M(R) = 1 - \max_R P_R$   is the parent loss;
$M(R_1) = 1 - \max_R P_{1R}$   is the first child loss;
$M(R_2) = 1 - \max_R P_{2R}$   is the second child loss
as long as $\max_R P_R = \max_R P_{1R} = \max_R P_{2R}$,

the weighted Gini loss will keep the same as parent.
$M(R) = q (1 - \max_R P_{1R}) + (1-q)(1 - \max_R P_{2k}) = M(R_1, R_2)$.

(b) : Random Forests.

   i. prove that:
$$Var(\hat{T}) = \rho\delta^2 + \frac{1-\rho}{B}\delta^2.$$

use the fact that:
$$Var(x+y) = Var(x) + Var(y) + 2Cov(x, y);$$

we have:
$$Var(\hat{T}) = Var\left(\frac{1}{B}\sum_{i=1}^{B} T_i(x)\right)$$

$$= \frac{1}{B^2} Var\left(\sum_{i=1}^{B} T_i\right)$$

$$= \frac{1}{B^2}\left(\sum_{i=1}^{B} Var(T_i) + 2\sum_{j=2}^{B}\sum_{i<j} Cov(T_i, T_j)\right)$$

$$= \frac{1}{B^2}\left(B\delta^2 + B(B-1)\cdot\rho\cdot\delta^2\right), \text{[choose 2 from B samples]}$$

$$= \delta^2\left(\frac{1}{B} + \frac{B-1}{B}\rho\right)$$

$$= \delta^2\left(\rho + \frac{1-\rho}{B}\right)$$

$$= \rho\delta^2 + \frac{1-\rho}{B}\delta^2, \quad \text{proved.}$$

   ii. Since $Var(\hat{T}) = \rho\delta^2 + \frac{1-\rho}{B}\delta^2$,

for the case of random forest it will have de-correlation effect, which means value $\rho$ above will be small, which means $Var(\hat{T})$ in random forest is going to be reduced. This makes bagging have higher variance.